



Margin conditions for vector quantization

Clément Levrard

► To cite this version:

| Clément Levrard. Margin conditions for vector quantization. 2014. hal-00877093v2

HAL Id: hal-00877093

<https://hal.archives-ouvertes.fr/hal-00877093v2>

Preprint submitted on 25 Apr 2014

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

MARGIN CONDITIONS FOR VECTOR QUANTIZATION

BY CLÉMENT LEVRARD

Université Paris Sud, UPMC and INRIA

Recent results in quantization theory show that the convergence rate for the mean-squared expected distortion of the empirical risk minimizer strategy, for any fixed probability distribution satisfying some regularity conditions, is $\mathcal{O}(1/n)$, where n is the sample size (see, e.g., [7] or [12]). However, the dependency of the average distortion on other parameters is not known.

This paper offers more general conditions, which may be thought of as margin conditions (see, e.g., [15]), under which a sharp upper bound on the expected distortion rate of the empirically optimal quantizer is derived. This upper bound is also proved to be sharp with respect to the dependency of the distortion on other natural parameters of the quantization issue.

1. Introduction. Quantization, also called lossy data compression in information theory, is the problem of replacing a probability distribution with an efficient and compact representation, that is a finite set of points. To be more precise, let P denote a probability distribution over \mathbb{R}^d and k a positive integer. A so-called k -quantizer Q is a map from \mathbb{R}^d to \mathbb{R}^d , whose image set is made of exactly k points, that is $|Q(\mathbb{R}^d)| = k$. For such a quantizer, every image point $c_i \in Q(\mathbb{R}^d)$ is called a code point, and the vector composed of the code points (c_1, \dots, c_k) is called a codebook. By considering the preimages of its code points, a quantizer Q partitions the Euclidean space \mathbb{R}^d into k groups, and assigns each group a representative. General references on the subject are to be found in [10], [9] and [13] among others.

The quantization theory was originally developed as a way to answer signal compression issues in the late 40's (see, e.g., [9]). However, unsupervised classification is also in the scope of its application. Isolating meaningful groups from a cloud of data is a topic of interest in many fields, from social science to biology. Classifying points into dissimilar groups of similar items is as more interesting as the amount of accessible data is large. In many cases data need to be preprocessed through a quantization algorithm in order to be exploited.

Keywords and phrases: localization, fast rates, margin conditions

If the distribution P has a finite second moment, the performance of a quantizer Q is measured by the risk, or distortion

$$R(Q) := P\|x - Q(x)\|^2,$$

where Pf means integration of the function f with respect to P . The choice of the Euclidean squared norm is convenient, since it takes advantages of the Euclidean space structure of \mathbb{R}^d . Nevertheless, it is worth pointing out that several authors deal with more general distortion functions. For further information on this topic, the interested reader is referred to [10] or [8].

In order to minimize the distortion introduced above, it is clear that only quantizers of the type $x \mapsto \arg \min_{c_1, \dots, c_k} \|x - c_i\|^2$ are to be considered. Such quantizers are called nearest-neighbor quantizers. With a slight abuse of notation, $R(\mathbf{c})$ will denote the risk of the nearest-neighbor quantizer associated with a codebook \mathbf{c} .

Provided that $P\|x\|^2 < \infty$, there exist optimal codebooks minimizing the risk R (see, e.g., Lemma 8 in [19] or Theorem 4.12 in [10]). The aim is to design a codebook $\hat{\mathbf{c}}_n$, according to a n -sample drawn from P , whose distortion is as close as possible to the optimal distortion $R(\mathbf{c}^*)$, where \mathbf{c}^* denotes an optimal codebook.

To solve this problem, most approaches to date attempt to implement the principle of empirical risk minimization in the vector quantization context. Let X_1, \dots, X_n denote an independent and identically distributed sample with distribution P . According to this principle, good code points can be found by searching for ones that minimize the empirical distortion over the training data, defined by

$$\hat{R}_n(\mathbf{c}) := \frac{1}{n} \sum_{i=1}^n \|X_i - Q(X_i)\|^2 = \frac{1}{n} \sum_{i=1}^n \min_{j=1, \dots, k} \|X_i - c_j\|^2.$$

If the training data represents the source well, then $\hat{\mathbf{c}}_n$ will hopefully also perform near optimally on the real source, that is $\ell(\hat{\mathbf{c}}_n, \mathbf{c}^*) = R(\hat{\mathbf{c}}_n) - R(\mathbf{c}^*) \approx 0$. The problem of quantifying how good empirically designed codebooks are, compared to the truly optimal ones, has been extensively studied, as for instance in [13].

It has been proved in [14] that $\mathbb{E}\ell(\hat{\mathbf{c}}_n, \mathbf{c}^*) = \mathcal{O}(1/\sqrt{n})$, provided that P has a finite second moment. However, this upper bound can be tightened whenever the source distribution satisfies additional assumptions.

For the special case of finitely supported distributions, it is shown in [2] that $\mathbb{E}\ell(\hat{\mathbf{c}}_n, \mathbf{c}^*) = \mathcal{O}(1/n)$. There are much more results in the case where P is assumed to have a density.

In fact, different sets of assumptions have been introduced in [2], [20] or [12], for the loss $\mathbb{E}\ell(\hat{\mathbf{c}}_n, \mathbf{c}^*)$ to decrease at the rate $\mathcal{O}(1/n)$ in the density case. As shown in [12], these different sets of assumptions turn out to be equivalent to a technical condition, similar to that used in [17] to derive fast rates of convergence in the statistical learning framework.

Thus, a question of interest is to know whether some margin type conditions can be derived for the source distribution to satisfy the technical condition mentioned above, as has been done in the statistical learning framework in [15].

Theorem 3.2 of [12] offers a partial answer, proving that a sufficient condition is that P is divided into k well separated areas. However, this condition is not fully satisfactory, since it consists in a bound on the density located at the $d - 1$ dimensional region between optimal code cells, whereas margin conditions in the statistical learning framework are bounds on the weight with respect to P of the ε -neighborhood of the critical value $1/2$ for the regression function.

Next, the scope of Theorem 3.2 of [12] is constrained to distributions with continuous densities, whereas margin conditions in [15] do not require regularity of the regression function.

This paper addresses both these issues, providing a condition which can clearly be thought of as a margin condition in the quantization framework, under which the loss $\mathbb{E}\ell(\hat{\mathbf{c}}_n, \mathbf{c}^*) = \mathcal{O}(1/n)$.

Moreover, some explicit oracle inequality is derived in this case, that is an upper bound of the form $\mathbb{E}\ell(\hat{\mathbf{c}}_n, \mathbf{c}^*) \leq C(k, d, P)/n$, where the dependency of $C(k, d, P)$ on its parameters is explicit, developing the technique used in [12] or [7]. It is worth pointing out that the parameters mentioned in this result, such as the smaller distance between two optimal code points, are rather natural from the quantization point of view.

In addition, this result allows to partially answer the problem mentioned in [1] about the minimax rates over distributions satisfying Pollard's condition. This rate has been proved in [1] to be $1/\sqrt{n}$, which is at first sight contradictory with the individual convergence rate of $1/n$ derived for every distribution in this case.

The paper is organized as follows. In Section 2 some notation and definition are introduced, as well as the so-called margin conditions. The main results are exposed in Section 3: firstly an oracle inequality on the loss is stated, along with a minimax result, then it is shown that Gaussian mixtures are in the scope of the margin conditions. Finally, proofs are gathered in Section 4, and the proofs of technical intermediate results are to be found in Section 5.

2. Notation and Definitions. Throughout the paper, for $M > 0$ and a in \mathbb{R}^d , $\mathcal{B}(a, M)$ will denote the closed Euclidean ball with center a and radius M . With a slight abuse of notation, P is said to be M -bounded if its support is included in $\mathcal{B}(0, M)$.

To frame the quantization issue as an empirical risk minimization issue, the following contrast function γ is introduced as

$$\gamma : \begin{cases} (\mathbb{R}^d)^k \times \mathbb{R}^d & \longrightarrow \mathbb{R} \\ (\mathbf{c}, x) & \longmapsto \min_{j=1, \dots, k} \|x - c_j\|^2, \end{cases}$$

where $\mathbf{c} = (c_1, \dots, c_k)$ denotes a codebook, that is a kd -dimensional vector. The risk $R(\mathbf{c})$ then takes the form $R(\mathbf{c}) = R(Q) = P\gamma(\mathbf{c}, \cdot)$, where we recall that Pf denotes the integration of the function f with respect to P . Similarly, the empirical risk $\hat{R}_n(\mathbf{c})$ can be defined as $\hat{R}_n(\mathbf{c}) = P_n\gamma(\mathbf{c}, \cdot)$, where P_n is the empirical distribution associated with X_1, \dots, X_n , in other words $P_n(A) = 1/n |\{i | X_i \in A\}|$, for every measurable subset $A \subset \mathbb{R}^d$.

It is worth pointing out that, if $P\|x\|^2 < \infty$, then there exist such minimizers $\hat{\mathbf{c}}_n$ and \mathbf{c}^* (see, e.g., Theorem 4.12 in [10]). In the sequel the set of minimizers of the risk $R(\cdot)$ will be denoted by \mathcal{M} .

Let c_1, \dots, c_k be a sequence of code points. A central role is played by the set of points which are closer to c_i than to any other c_j 's. To be more precise, the Voronoi cell, or quantization cell associated with c_i is the closed set defined by

$$V_i(\mathbf{c}) = \left\{ x \in \mathbb{R}^d \mid \forall j \neq i \quad \|x - c_i\| \leq \|x - c_j\| \right\}.$$

It may be noted that $(V_1(\mathbf{c}), \dots, V_k(\mathbf{c}))$ does not form a partition of \mathbb{R}^d , since $V_i(\mathbf{c}) \cap V_j(\mathbf{c})$ may be non empty. To address this issue, a Voronoi partition associated with \mathbf{c} is defined as a sequence of subsets $(W_1(\mathbf{c}), \dots, W_k(\mathbf{c}))$ which forms a partition of \mathbb{R}^d , and such that for every $i = 1, \dots, k$,

$$\bar{W}_i(\mathbf{c}) = V_i(\mathbf{c}),$$

where $\bar{W}_i(\mathbf{c})$ denotes the closure of the subset $W_i(\mathbf{c})$. The open Voronoi cell is defined the same way by

$$\overset{o}{V}_i(\mathbf{c}) = \left\{ x \in \mathbb{R}^d \mid \forall j \neq i \quad \|x - c_i\| < \|x - c_j\| \right\}.$$

Given a Voronoi partition $W(\mathbf{c}) = (W_1(\mathbf{c}), \dots, W_k(\mathbf{c}))$, the following inclusion holds, for i in $\{1, \dots, k\}$,

$$\overset{o}{V}_i(\mathbf{c}) \subset W_i(\mathbf{c}) \subset V_i(\mathbf{c}),$$

and the risk $R(\mathbf{c})$ takes the form

$$R(\mathbf{c}) = \sum_{i=1}^k P(\|x - c_i\|^2 1_{W_i(\mathbf{c})}(x)),$$

where 1_A denotes the indicator function associated with A . In the case where (W_1, \dots, W_k) are fixed subsets such that $P(W_i) \neq 0$, for every $i = 1, \dots, k$, it is clear that

$$P(\|x - c_i\|^2 1_{W_i(\mathbf{c})}(x)) \geq P(\|x - \eta_i\|^2 1_{W_i(\mathbf{c})}(x)),$$

where η_i denotes the conditional expectation of P over the subset $W_i(\mathbf{c})$, that is

$$\eta_i = \frac{P(x 1_{W_i(\mathbf{c})}(x))}{P(W_i(\mathbf{c}))}.$$

Moreover, it is proved in Theorem 4.1 of [10] that, for every Voronoi partition $W(\mathbf{c}^*)$ associated with an optimal codebook \mathbf{c}^* , and every $i = 1, \dots, k$, $P(W_i(\mathbf{c}^*)) \neq 0$. Consequently, any optimal codebook satisfies the so-called centroid condition (see, e.g., Section 6.2 of [9]), that is

$$\mathbf{c}_i^* = \frac{P(x 1_{W_i(\mathbf{c}^*)}(x))}{P(W_i(\mathbf{c}^*))}.$$

As a remark, the centroid condition ensures that, for every \mathbf{c}^* in \mathcal{M} and $i \neq j$,

$$P(V_i(\mathbf{c}^*) \cap V_j(\mathbf{c}^*)) = P\left(\left\{x \in \mathbb{R}^d \mid \|x - c_i^*\| = \|x - c_j^*\|\right\}\right) = 0.$$

A proof of this statement can be found in Theorem 4.2 of [10]. According to this remark, it is clear that, for every optimal Voronoi partition $(W_1(\mathbf{c}^*), \dots, W_k(\mathbf{c}^*))$,

$$(1) \quad \begin{cases} P(W_i(\mathbf{c}^*)) &= P(V_i(\mathbf{c}^*)) \\ P_n(W_i(\mathbf{c}^*)) &\underset{a.s.}{=} P_n(V_i(\mathbf{c}^*)). \end{cases}$$

The following quantities are of importance in the bounds exposed in Section 3.1:

$$(2) \quad \begin{cases} B &= \min_{\mathbf{c}^* \in \mathcal{M}, i \neq j} \|c_i^* - c_j^*\| \\ p_{min} &= \min_{\mathbf{c}^* \in \mathcal{M}, i=1, \dots, k} P(V_i(\mathbf{c}^*)). \end{cases}$$

The role of the boundaries between optimal Voronoi cells may be compared to the role played by the critical value $1/2$ for the regression function

in the statistical learning framework. To draw this comparison, the following set is introduced, for any $\mathbf{c}^* \in \mathcal{M}$,

$$N(\mathbf{c}^*) = \bigcup_{i \neq j} V_i(\mathbf{c}^*) \cap V_j(\mathbf{c}^*).$$

Next, the critical region N^* is defined as

$$N^* = \bigcup_{\mathbf{c}^* \in \mathcal{M}} N(\mathbf{c}^*).$$

This region seems to be of importance when considering the conditions under which the empirical risk minimization strategy for the quantization issue achieves faster rates of convergence, as exposed in [12]. However, to fully draw the comparison between the margin conditions for the statistical learning issue (see, e.g., [15]) and quantization, the neighborhood of this region has to be introduced. For this purpose the t -neighborhood of the critical region is defined as

$$N_t^* = \left\{ x \in \mathbb{R}^d \mid d(x, N^*) \leq t \right\}.$$

Intuitively, if $P(N_t^*)$ is small enough, then the source distribution P is concentrated around its optimal codebook, and may be thought of as a slight modification of the probability distribution with finite support made of an optimal codebook \mathbf{c}^* . To be more precise, let introduce the following key assumption:

DEFINITION 2.1 (Margin condition). *Denote by $p(t) = P(N^*(t))$. Then P satisfies a margin condition with radius r_0 if and only if*

- i) P is bounded by M ,*
- ii) \mathcal{M} is finite,*
- iii) For all $0 \leq t \leq r_0$,*

$$(3) \quad p(t) \leq \frac{Bp_{\min}}{128M^2}t.$$

Contrary to the conditions required in [15] in the framework of supervised classification, the margin condition introduced here only requires a local control of the weight of the neighborhood of the critical region. It is quite obvious that a global margin condition of the type $p(t) \leq Bp_{\min}/128M^2t$ for every $t > 0$ implies the condition defined above. However, requiring only a local control of the weight function $p(t)$ enlarges the scope of our results,

since it allows to deal with non continuous probability distributions. This point is illustrated in the following example:

Example 1: Assume that there exists $r > 0$ such that $p(x) = 0$ if $x \leq r$ (for instance if P is supported on k points). Then P satisfies a margin condition with radius r .

It is also worth pointing out that the condition mentioned in [15] requires a control of the weight of the neighborhood of the critical value $1/2$ with a polynomial function with degree larger than 1. In the quantization framework, the special role played by the exponent 1 leads to only consider linear controls of the weight function. This point is explained by the following example:

Example 2: Assume that P is bounded by M , and that there exists $Q > 0$ and $q > 1$ such that $p(x) \leq Qx^q$. Then P satisfies (3), with

$$r_0 = \frac{B}{4\sqrt{2}M} \left(\frac{p_{\min} B}{16\sqrt{2}MQ} \right)^{1/(q-1)}.$$

In the case where P has a density, the condition (3) can be thought of as a generalization of the condition mentioned in Theorem 3.2 of [12], which requires the density of the distribution to be small enough over the critical region. In fact, provided that P has a continuous density, a uniform bound on the density over the critical region provides a local control of the weight function with a polynomial function of degree 1. This idea is developed in the following example:

Example 3(Continuous densities): Assume that P has a continuous density f , is bounded by M , and that \mathcal{M} is finite. Moreover, assume that

$$(4) \quad \int_{N^*} f(u) du < \frac{B p_{\min}}{128M^2}.$$

Then, by considering the derivative at 0 of the map $t \mapsto p(t)$, there exists $r_0 > 0$ such that P satisfies a margin condition with radius $r_0 > 0$. It can easily be deduced from (4) that an uniform bound on the density located at the critical region can provide a sufficient condition for a distribution P to satisfy a margin condition. Such a result has to be compared to Theorem 3.2 of [12], where it was required that

$$\|f|_{N^*}\|_{\infty} \leq \frac{\Gamma(\frac{d}{2}) B}{2^{d+5} M^{d+1} \pi^{d/2} p_{\min}},$$

where Γ denotes the Gamma function.

Another interesting parameter of the quantization issue is the following separation factor, which quantifies the difference between optimal codebooks and local minimizers of the risk.

DEFINITION 2.2. Denote by $\tilde{\mathcal{M}}$ the set of local minimizers of the map $\mathbf{c} \mapsto P\gamma(\mathbf{c}, \cdot)$. Then P is said to be ε -separated if

$$(5) \quad \inf_{\mathbf{c} \in \tilde{\mathcal{M}} \cap \mathcal{M}^c} \ell(\mathbf{c}, \mathbf{c}^*) = \varepsilon.$$

It may be noticed that local minimizers of the risk function satisfy the centroid condition. Whenever P has a density and $P\|x\|^2 < \infty$, it can be proved that the set of minimizers of R coincides with the set of codebooks satisfying the centroid condition, also called stationary points (see, e.g., Lemma A of [20]). However, this result cannot be extended to non continuous distributions, as proved in Example 4.11 of [10].

The main results of the present paper are based on the following proposition, which connects the margin condition stated in Definition 2.1 to the condition introduced in Theorem 2 of [2].

PROPOSITION 2.1. Assume that P satisfies a margin condition with radius r_0 , and is ε -separated. Then, for every codebook \mathbf{c} in $\mathcal{B}(0, M)^k$,

$$\|\mathbf{c} - \mathbf{c}^*(\mathbf{c})\|^2 \leq \kappa_0 \ell(\mathbf{c}, \mathbf{c}^*),$$

where $\mathbf{c}^*(\mathbf{c}) \in \arg \min_{\mathbf{c}^* \in \mathcal{M}} \|\mathbf{c} - \mathbf{c}^*\|$, and $\kappa_0 = 4kM^2 \left(\frac{1}{\varepsilon} \vee \frac{64M^2}{p_{\min} B^2 r_0^2} \right)$.

As mentioned in [7] or [12], the connection between the loss and the Euclidean squared distance can be thought of as a technical margin condition. It is worth pointing out that the dependency of κ_0 on different parameters of the quantization issue is explicit. This point allows us to derive explicit upper bounds on the excess risk in the following section.

3. Results.

3.1. *Risk bound.* The main result of this paper is the following:

THEOREM 3.1. Assume that P satisfies a margin condition with radius r_0 , and is ε -separated. Let κ_0 be defined as

$$\kappa_0 = 4kM^2 \left(\frac{1}{\varepsilon} \vee \frac{64M^2}{p_{\min} B^2 r_0^2} \right).$$

If $\hat{\mathbf{c}}_n$ is an empirical risk minimizer, then, with probability larger than $1 - 2e^{-x}$,

$$(6) \quad \ell(\hat{\mathbf{c}}_n, \mathbf{c}^*) \leq C_0 \kappa_0 \frac{|\mathcal{M}|^2 R(\mathbf{c}^*)}{n} + \kappa_0^3 \frac{C_1}{n^2} + \kappa_0 \frac{C_2}{n} x + \frac{C_3}{n},$$

where C_0 is an absolute constant, C_1 is a combination of square roots of polynomial functions in k , $\log(k)$, d , B and M , C_2 is polynomial in k and M , C_3 is polynomial in M and \sqrt{k} .

Moreover, under the same conditions, with probability larger than $1 - 2e^{-x}$,

$$(7) \quad \ell(\hat{\mathbf{c}}_n, \mathbf{c}^*) \leq C'_0 \kappa_0 \frac{M^2 k d (\log(4|\mathcal{M}|\sqrt{k d}) + 1)}{n} + \kappa_0 \frac{144 M^2}{n} x + \frac{64 M^2}{n} x,$$

where C'_0 is an absolute constant.

This result is in line with Theorem 3.1 in [12] or Theorem 1 in [7], concerning the dependency on the sample size n of the loss $\ell(\hat{\mathbf{c}}_n, \mathbf{c}^*)$. The main advance lies in the dependency on other parameters of the loss of $\hat{\mathbf{c}}_n$, which provides a non-asymptotic bound for the excess risk.

In fact, (7) derives from chaining arguments such as one used in [12] or [7], and involves a classical dimension term of kd . When considering (6), it seems that this kd term disappears from the dominant term of the upper bound. This suggests that the dimension of the Euclidean space in the finite-dimensional case plays a minor role, as pointed out in Theorem 2.1 in [5]. An open question is to know whether such fast rates bounds can be derived in the infinite dimensional case.

However, (6) may be thought of as a semi-asymptotic bound, since it involves a dominant term and a residual term with respect to the sample size n . Although the dependency on other parameters of the dominant term is sharper in (6) than in (7), the residual term C_1/n^2 in (6) still involves the dimension d . Consequently, (6) only guarantees that $\mathbb{E}\ell(\hat{\mathbf{c}}_n, \mathbf{c}^*)$ can be bounded from above with a dimension-free term when n grows to infinity.

Another interesting point is that Theorem 3.1 does not require P to have a density, contrary to the requirements of previous results in [12] or [7]. This remark makes the link between bounds obtained for point-wise distributions in [2] and bounds for distributions with densities as in [2] or [12].

It is also worth mentioning that the dependency in ε surprisingly turns out to be sharp, as will be shown in Proposition 3.1. In fact, tuning this separation factor is the core of the demonstration of the minimax results in [3] or [1].

3.2. Minimax lower bound. Theorem 1 in [3] ensures that the minimax convergence rate over the distributions bounded by M of any empirically designed codebook can be bounded from below by $\mathcal{O}(1/\sqrt{n})$. A question of interest is to know whether this lower bound can be refined when considering

only distributions satisfying some fast-convergence condition. A partial answer is given by Corollary 2 in [1], where it is proved that the minimax rate over distributions with continuous densities with individual convergence rate of $\mathcal{O}(1/n)$ for the empirical risk minimizer is still $\mathcal{O}(1/\sqrt{n})$. However, since no non-asymptotic upper bound has been provided for these distributions, to understand which parameter is varying in this minimax result remains a hard issue.

Consequently, this subsection is devoted to obtaining a minimax lower bound on the excess risk over the set of distributions satisfying the margin condition defined in Definition 2.1, in which some parameters are fixed. Throughout this subsection, $\hat{\mathbf{c}}_n$ will denote an empirically designed codebook, that is a map from $(\mathbb{R}^d)^n$ to $(\mathbb{R}^d)^k$. Let k be an integer such that $k \geq 3$, and $M > 0$. For simplicity, k is assumed to be divisible by 3. Let us introduce the following quantities:

$$\begin{cases} m &= \frac{2k}{3} \\ \Delta &= \frac{15M}{96m^{1/d}}. \end{cases}$$

To focus on the dependency on the separation factor ε , the quantities involved in Definition 2.1 are fixed as:

$$(8) \quad \begin{cases} B &= \Delta \\ r_0 &= \frac{7\Delta}{16} \\ p_{\min} &\geq \frac{1}{2k}. \end{cases}$$

Denote by $\mathcal{D}(\varepsilon)$ the set of probability distributions which are ε -separated, and which satisfies a margin condition with parameters defined in (8). The minimax result is the following:

PROPOSITION 3.1. *Assume that $k \geq 3$. Then, for any empirically designed codebook,*

$$\mathbb{E} \sup_{P \in \mathcal{D}(c_1/\sqrt{n})} \ell(\hat{\mathbf{c}}_n, \mathbf{c}^*) \geq c_0 M^2 \frac{\sqrt{k^{1-\frac{4}{d}}}}{\sqrt{n}},$$

where c_0 is an absolute constant, and

$$c_1 = \frac{(15M)^2}{4(96m^{\frac{1}{4}+\frac{1}{d}})^2}.$$

Proposition 3.1 can be thought of as an extension of Theorem 1 in [3]. This minimax lower bound has to be compared to the upper risk bound

obtained in Theorem 3.1 for the empirical risk minimizer $\hat{\mathbf{c}}_n$ over the set of distributions $\mathcal{D}(c_1/\sqrt{n})$. To be more precise, Theorem 3.1 ensures that, provided that n is large enough,

$$\mathbb{E} \sup_{P \in \mathcal{D}(c_1/\sqrt{n})} \leq \frac{g(k, d, M)}{\sqrt{n}},$$

where $g(k, d, M)$ depends only on k , d and M . In other words, the dependency of the upper bounds stated in Theorem 3.1 on ε turns out to be sharp whenever $\varepsilon \sim n^{-\frac{1}{2}}$. Unfortunately, Proposition 3.1 can not be easily extended to the case where $\varepsilon \sim n^{-\alpha}$, with $0 < \alpha < 1/2$. Consequently an open question is whether the upper bounds stated in Theorem 3.1 remains accurate with respect to ε in this case.

3.3. Quasi-Gaussian mixture example. The aim of this subsection is to illustrate the results offered in Section 3 with Gaussian mixtures in dimension $d = 2$. The Gaussian mixture model is a typical and well-defined clustering example. However we will not deal with the clustering issue but rather with its theoretical background.

In general, a Gaussian mixture distribution \tilde{P} is defined by its density

$$\tilde{f}(x) = \sum_{i=1}^{\tilde{k}} \frac{\theta_i}{2\pi\sqrt{|\Sigma_i|}} e^{-\frac{1}{2}(x-m_i)^t \Sigma_i^{-1} (x-m_i)},$$

where \tilde{k} denotes the number of component of the mixture, and the θ_i 's denote the weights of the mixture, which satisfy $\sum_{i=1}^{\tilde{k}} \theta_i = 1$. Moreover, the m_i 's denote the means of the mixture, so that $m_i \in \mathbb{R}^2$, and the Σ_i 's are the 2×2 variance matrices of the components.

We restrict ourselves to the case where the number of components \tilde{k} is known, and match the size k of the codebooks. To ease the calculation, we make the additional assumption that every component has the same diagonal variance matrix $\Sigma_i = \sigma^2 I_2$. Note that a similar result to Proposition 3.2 can be derived for distributions with different variance matrices Σ_i , at the cost of more computing.

Since the support of a Gaussian random variable is not bounded, we define the “quasi-Gaussian” mixture model as follows, truncating each Gaussian component. Let the density f of the distribution P be defined by

$$f(x) = \sum_{i=1}^k \frac{\theta_i}{2\pi\sigma^2 N_i} e^{-\frac{\|x-m_i\|^2}{2\sigma^2}} 1_{\mathcal{B}(0, M)},$$

where N_i denotes a normalization constant for each Gaussian variable.

To ensure this model to be close to the Gaussian mixture model, we assume that there exists a constant $\varepsilon \in [0, 1]$ such that, for $i = 1, \dots, k$, $N_i \geq 1 - \varepsilon$.

Denote by $\tilde{B} = \inf_{i \neq j} \|m_i - m_j\|$ the smallest possible distance between two different means of the mixture. To avoid boundary issues we assume that, for all $i = 1, \dots, k$, $\mathcal{B}(m_i, \tilde{B}/3) \subset \mathcal{B}(0, M)$.

It is worth noticing that the two assumptions $N_i \geq 1 - \varepsilon$ and $\mathcal{B}(m_i, \tilde{B}/3) \subset \mathcal{B}(0, M)$ can easily be satisfied as soon as M is chosen large enough. For such a model, Proposition 3.2 offers a sufficient condition for P to satisfy a margin condition.

PROPOSITION 3.2. *Let $\theta_{\min} = \min_{i=1, \dots, k} \theta_i$, and $\theta_{\max} = \max_{i=1, \dots, k} \theta_i$. Assume that*

$$(9) \quad \frac{\theta_{\min}}{\theta_{\max}} \geq \max \left(\frac{2048k\sigma^2}{(1 - \varepsilon)\tilde{B}^2(1 - e^{-\tilde{B}^2/2048\sigma^2})}, \frac{2048k^2M^3}{(1 - \varepsilon)7\sigma^2\tilde{B}(e^{\tilde{B}^2/32\sigma^2} - 1)} \right).$$

Then P satisfies a margin condition with radius $\frac{\tilde{B}}{8}$.

The condition (9) can be decomposed as follows. If

$$\frac{\theta_{\min}}{\theta_{\max}} \geq \frac{2048k\sigma^2}{(1 - \varepsilon)\tilde{B}^2(1 - e^{-\tilde{B}^2/2048\sigma^2})},$$

then the optimal codebook \mathbf{c}^* is close to the vector of means of the mixture $\mathbf{m} = (m_1, \dots, m_k)$. Therefore, it is possible to locate the critical region associated with the optimal codebook \mathbf{c}^* , and to derive an upper bound on the weight function defined in Definition 2.1. This leads to the second term of the maximum in (9).

This condition can be interpreted as a condition on the polarization of the mixture. A favorable case for vector quantization seems to be when the poles of the mixtures are well separated, which is equivalent to σ is small compared to \tilde{B} , when considering Gaussian mixtures. Proposition 3.2 gives details on how σ has to be small compared to \tilde{B} , in order to satisfy the requirements of Proposition 2.1. This ensures that the loss $\ell(\hat{\mathbf{c}}_n, \mathbf{c}^*)$ reaches an improved convergence rate of $1/n$.

It may be noticed that Proposition 3.2 offers almost the same condition than Proposition 4.2 in [12]. In fact, since the Gaussian mixture distributions have a continuous density, making use of (4) in Example 3 ensures that the

margin condition for Gaussian mixtures is equivalent to a bound on the density over the critical region.

It is important to note that this result is valid when k is known and match exactly the number of components of the mixture. When the number of code points k is different from the number of components \tilde{k} of the mixture, we have no general idea of where the optimal code points can be located.

Moreover, suppose that there exists only one optimal codebook \mathbf{c}^* , up to reindexing, and that we are able to locate this optimal codebook \mathbf{c}^* . As mentioned in Proposition 2.1, the key quantity is in fact $B = \inf_{i \neq j} \|c_i^* - c_j^*\|$. In the case where $\tilde{k} \neq k$, there is no simple relation between \tilde{B} and B . Consequently, a condition like in Proposition 3.2 could not involve the natural parameter of the mixture \tilde{B} .

It is also worth pointing out that there exist cases where the set of optimal codebooks is not finite. For example, assume that P is a truncated rotationally symmetric Gaussian distribution, and $k = 2$. Since every rotation of an optimal codebook leads to another optimal codebook, there exists an infinite set of optimal codebooks. Since, in this case, $N^* = \mathcal{B}(0, M)$, obviously P does not satisfy a margin condition.

4. Proofs.

4.1. *Proof of Proposition 2.1.* The proof of Proposition 2.1 is based on the following lemma.

LEMMA 4.1. *Let $x \in V_i(\mathbf{c}^*) \cap V_j(\mathbf{c})$, for $i \neq j$. Then*

$$(10) \quad \left| \left\langle x - \frac{c_i + c_j}{2}, c_i - c_j \right\rangle \right| \leq 4\sqrt{2}M \|\mathbf{c} - \mathbf{c}^*\|,$$

$$(11) \quad d(x, \partial V_i(\mathbf{c}^*)) \leq \frac{4\sqrt{2}M}{B} \|\mathbf{c} - \mathbf{c}^*\|.$$

The two statements of Lemma 4.1 emphasize the fact that, provided that \mathbf{c} and \mathbf{c}^* are quite similar, the areas on which the label may differ with respect to \mathbf{c} and \mathbf{c}^* should be close to the boundary of Voronoi diagrams. This idea is mentioned in the proof of Corollary 1 in [2]. Nevertheless we provide here a simpler proof.

PROOF OF LEMMA 4.1. Let $x \in V_i(\mathbf{c}^*) \cap V_j(\mathbf{c})$, then $\|x - c_j\|^2 \leq \|x - c_i\|^2$, which leads to $\left\langle c_i - c_j, x - \frac{c_i + c_j}{2} \right\rangle \leq 0$. Since $\|x - c_i^*\| \leq \|x - c_j^*\|$, we may write

$$\|x - c_i\| \leq \|x - c_j\| + \|c_i - c_i^*\| + \|c_j - c_j^*\|.$$

Taking square on both sides leads to

$$\begin{aligned}
\|x - c_i\|^2 - \|x - c_j\|^2 &\leq 2\|x - c_j\|(\|c_i - c_i^*\| + \|c_j - c_j^*\|) \\
&\quad + (\|c_i - c_i^*\| + \|c_j - c_j^*\|)^2 \\
&\leq 8M(\|c_i - c_i^*\| + \|c_j - c_j^*\|) \\
&\leq 8\sqrt{2}M\|\mathbf{c} - \mathbf{c}^*\|.
\end{aligned}$$

Since $\|x - c_i\|^2 - \|x - c_j\|^2 = 2\left\langle x - \frac{c_i + c_j}{2}, c_i - c_j \right\rangle$, (10) is proved.

To prove (11), remark that, since $x \in V_i(\mathbf{c}^*)$, $d(x, V_i(\mathbf{c}^*)) \leq d(x, h_{i,j}^*)$, where $h_{i,j}^*$ is the hyperplane defined by $\left\{x \in \mathcal{B}(0, M) \mid \|x - c_i^*\| = \|x - c_j^*\|\right\}$. Using quite simple geometric arguments, we deduce that

$$d(x, h_{i,j}^*) = \left| \left\langle x - \frac{c_i^* + c_j^*}{2}, \frac{c_i^* - c_j^*}{\|c_i^* - c_j^*\|} \right\rangle \right|.$$

The same arguments as in the proof of (10) guarantee that

$$\begin{aligned}
\left| \left\langle x - \frac{c_i^* + c_j^*}{2}, \frac{c_i^* - c_j^*}{\|c_i^* - c_j^*\|} \right\rangle \right| &= \left\langle x - \frac{c_i^* + c_j^*}{2}, \frac{c_i^* - c_j^*}{\|c_i^* - c_j^*\|} \right\rangle \\
&\leq \frac{4\sqrt{2}M}{B}\|\mathbf{c} - \mathbf{c}^*\|.
\end{aligned}$$

□

Equipped with Lemma 4.1, we are in a position to prove Proposition 2.1. Let $\mathbf{c} \in \mathcal{M}$, and $(W_1(\mathbf{c}), \dots, W_k(\mathbf{c}))$ be a Voronoi partition associated to \mathbf{c} , as defined in Section 2. Then $\ell(\mathbf{c}, \mathbf{c}^*)$ can be decomposed as follows:

$$\begin{aligned}
P\gamma(\mathbf{c}, \cdot) &= \sum_{i=1}^k P(\|x - c_i\|^2 1_{W_i(\mathbf{c})}) \\
&= \sum_{i=1}^k P(\|x - c_i\|^2 1_{V_i(\mathbf{c}^*)}) + \sum_{i=1}^k P(\|x - c_i\|^2 (1_{W_i(\mathbf{c})} - 1_{V_i(\mathbf{c}^*)})).
\end{aligned}$$

Since, for all $i = 1, \dots, k$, $P(x 1_{V_i(\mathbf{c}^*)}(x)) = P(V_i(\mathbf{c}^*))c_i^*$ (centroid condition), we may write

$$P(\|x - c_i\|^2 1_{V_i(\mathbf{c}^*)}) = P(V_i(\mathbf{c}^*))\|c_i - c_i^*\|^2 + P(\|x - c_i^*\|^2 1_{V_i(\mathbf{c}^*)}),$$

from which we deduce

$$P\gamma(\mathbf{c}, \cdot) = P\gamma(\mathbf{c}^*, \cdot) + \sum_{i=1}^k P(V_i(\mathbf{c}^*)) \|c_i - c_i^*\|^2 + \sum_{i=1}^k P(\|x - c_i\|^2 (1_{W_i(\mathbf{c})} - 1_{V_i(\mathbf{c}^*)})),$$

which leads to

$$\ell(\mathbf{c}, \mathbf{c}^*) \geq p_{\min} \|\mathbf{c} - \mathbf{c}^*\|^2 + \sum_{i=1}^k \sum_{j \neq i} P \left((\|x - c_j\|^2 - \|x - c_i\|^2) 1_{V_i(\mathbf{c}^*) \cap W_j(\mathbf{c})} \right).$$

Since $x \in W_j(\mathbf{c}) \subset V_j(\mathbf{c})$, $\|x - c_j\|^2 - \|x - c_i\|^2 \leq 0$. Thus it remains to bound from above

$$\sum_{i=1}^k \sum_{j \neq i} P \left((\|x - c_i\|^2 - \|x - c_j\|^2) 1_{V_i(\mathbf{c}^*) \cap W_j(\mathbf{c})} \right).$$

Noticing that

$$\|x - c_i\|^2 - \|x - c_j\|^2 = 2 \left\langle c_j - c_i, x_{i,j} - \frac{c_i + c_j}{2} \right\rangle,$$

and using Lemma 4.1, we get

$$\sum_{i=1}^k P(\|x - c_i\|^2 (1_{W_i(\mathbf{c})} - 1_{V_i(\mathbf{c}^*)})) \geq -8\sqrt{2}M \|\mathbf{c} - \mathbf{c}^*\| p \left(\frac{4\sqrt{2}M}{B} \|\mathbf{c} - \mathbf{c}^*\| \right).$$

Consequently, if P satisfies (3), then, if $\|\mathbf{c} - \mathbf{c}^*\| \leq \frac{Br_0}{4\sqrt{2}M}$,

$$\ell(\mathbf{c}, \mathbf{c}^*) \geq \frac{p_{\min}}{2} \|\mathbf{c} - \mathbf{c}^*\|^2.$$

Now turn to the case where $\|\mathbf{c} - \mathbf{c}^*(\mathbf{c})\| \geq \frac{Br_0}{4\sqrt{2}M}$. Since the support of P is included on $\mathcal{B}(0, M)$, the function $\mathbf{c} \mapsto P\gamma(\mathbf{c}, \cdot)$ is continuous, its minimum on $(\mathbb{R}^d)^k \cap (\bigcup_{\mathbf{c}^* \in \mathcal{M}} \mathcal{B}(0, M))^c$ is attained. Such a minimizer is a local minimizer, or is at the boundary $\|\mathbf{c} - \mathbf{c}^*(\mathbf{c})\| = \frac{Br_0}{4\sqrt{2}M}$. Hence we deduce

$$\begin{aligned} \ell(\mathbf{c}, \mathbf{c}^*) &\geq \varepsilon \wedge \frac{p_{\min} Br_0^2}{64M^2} \\ &\geq \left(\varepsilon \wedge \frac{p_{\min} Br_0^2}{64M^2} \right) \frac{\|\mathbf{c} - \mathbf{c}^*\|^2}{4kM^2}. \end{aligned}$$

This proves Proposition 2.1

4.2. *Proof of Theorem 3.1.* Throughout this subsection P is assumed to satisfy a margin condition with radius r_0 , and to be ε -separated. A non decreasing map $\Phi : \mathbb{R} \rightarrow \mathbb{R}^+$ is called sub- α if $x \mapsto \frac{\Phi(x)}{x^\alpha}$ is non increasing.

The following localization theorem, derived from Theorem 6.1 in [6], is the main argument of our proof.

THEOREM 4.1. *Let \mathcal{F} be a class of bounded measurable functions such that there exist $b > 0$ and $\omega : \mathcal{F} \rightarrow \mathbb{R}^+$ satisfying*

- (i) $\forall f \in \mathcal{F} \quad \|f\|_\infty \leq b,$
- (ii) $\forall f \in \mathcal{F} \quad \text{Var}(f) \leq \omega(f).$

Let K be a positive constant, Φ a sub- α function, $\alpha \in [1/2, 1]$. Then there exists a constant $C(\alpha)$ such that, if D is a constant satisfying $D \leq 6KC(\alpha)$, and r^ is the unique solution of the equation $\Phi(r) = r/D$, the following holds. Assume that*

$$\forall r \geq r^* \quad \mathbb{E} \left(\sup_{\omega(f) \leq r} |(P - P_n)f| \right) \leq \Phi(r).$$

Then, for all $x > 0$, with probability larger than $1 - e^{-x}$,

$$\forall f \in \mathcal{F} \quad Pf - P_n f \leq K^{-1} \left(\omega(f) + \left(\frac{6KC(\alpha)}{D} \right)^{\frac{1}{1-\alpha}} r^* + \frac{(9K^2 + 16Kb)x}{4n} \right).$$

A proof of Theorem 4.1 is given in Section 5.3 of [12]. Notice that an explicit calculation of $C(\alpha)$ is given by $C(\alpha) = \inf_{x>1} \left(1 + x^\alpha \left(\frac{1}{2} + \frac{1}{x^{1-\alpha}-1} \right) \right)$.

4.2.1. *Proof of (6).* The proof of (6) follows from the combination of Proposition 2.1 and a direct application of Theorem 4.1. To be more precise, let \mathcal{F}_1 denote the set

$$\mathcal{F}_1 = \left\{ \gamma(\mathbf{c}, \cdot) - \gamma(\mathbf{c}^*(\mathbf{c}), \cdot) \mid \mathbf{c} \in \mathcal{B}(0, M)^k \right\}.$$

Since, for all $i \in \{1, \dots, k\}$,

$$\left| \|x - c_i\|^2 - \|x - c_i^*(\mathbf{c})\|^2 \right| \leq 4M \|c_i - c_i^*(\mathbf{c})\|,$$

it follows that, for every $f \in \mathcal{F}_1$,

$$\begin{cases} \|f\|_\infty & \leq 8M^2 \\ \text{Var}_P(f) & \leq 16M^2 \|\mathbf{c} - \mathbf{c}^*(\mathbf{c})\|^2. \end{cases}$$

Define $\omega_1(f) = 16M^2 \|\mathbf{c} - \mathbf{c}^*(\mathbf{c})\|^2$. It remains to bound from above the complexity term. This is done in the following proposition, derived from the proof of Theorem 1 in [7].

PROPOSITION 4.1. *One has*

$$(12) \quad \mathbb{E} \sup_{f \in \mathcal{F}_1, \omega_1(f) \leq \delta} |(P - P_n)f| \leq \frac{(2\sqrt{2} + 64)\sqrt{kd}}{\sqrt{n}} \left(\sqrt{\log(4|\mathcal{M}|\sqrt{kd})} + 1 \right) \sqrt{\delta}.$$

The proof of Proposition 4.1 derives from classical chaining arguments, and is given in Section 5.1. Let Φ_1 be defined as the right-hand side of (12). Observing that $\Phi_1(\delta)$ takes the form $\Phi_1(\delta) = \Xi_1 \sqrt{\delta/n}$, the solution δ_1^* of the equation $\Phi_1(\delta) = \delta/D$ may be written, for any $D > 0$,

$$\delta_1^* = \frac{D^2 \Xi_1^2}{n}.$$

Let $K > 0$ and choose $D = 6KC(1/2)$. Applying Theorem 4.1 to \mathcal{F}_1 leads to, with probability larger than $1 - e^{-x}$,

$$\begin{aligned} (P - P_n)(\gamma(\mathbf{c}, \cdot) - \gamma(\mathbf{c}^*(\mathbf{c}), \cdot)) &\leq K^{-1} 16M^2 \|\mathbf{c} - \mathbf{c}^*(\mathbf{c})\|^2 \\ &\quad + \frac{36KC(1/2)^2 \Xi_1^2}{n} + \frac{9K + 128M^2}{4n} x. \end{aligned}$$

Introducing the inequality $\kappa_0 \ell(\mathbf{c}, \mathbf{c}^*) \geq \|\mathbf{c} - \mathbf{c}^*(\mathbf{c})\|^2$ provided by Proposition 2.1, choosing $K = 32M^2 \kappa_0$ and taking into account that $C(1/2) \leq 4$ leads to (6).

4.2.2. *Proof of (7).* The proof of (7) also relies on an application of Theorem 4.1. Let the loss of $\hat{\mathbf{c}}_n$ be decomposed as follows,

$$\begin{aligned} (13) \quad P(\gamma(\hat{\mathbf{c}}_n) - \gamma(\mathbf{c}^*)) &\leq (P - P_n)(\gamma(\hat{\mathbf{c}}_n) - \gamma(\mathbf{c}^*)) \\ &\leq (P - P_n) \langle \hat{\mathbf{c}}_n - \mathbf{c}^*(\hat{\mathbf{c}}_n), \Delta(\mathbf{c}^*(\hat{\mathbf{c}}_n), \cdot) \rangle \\ &\quad + (P - P_n) \|\hat{\mathbf{c}}_n - \mathbf{c}^*(\hat{\mathbf{c}}_n)\| R(\hat{\mathbf{c}}_n, \mathbf{c}^*(\hat{\mathbf{c}}_n), \cdot), \end{aligned}$$

where

$$\Delta(\mathbf{c}^*, x) = -2((x - c_1^*)1_{V_1(\mathbf{c}^*)}, \dots, (x - c_k^*)1_{V_k(\mathbf{c}^*)}),$$

and

$$\begin{aligned} R(\mathbf{c}, \mathbf{c}^*, x) &= \sum_{i,j=1,\dots,k} 1_{V_i(\mathbf{c}^*) \cap W_j(\mathbf{c})} \|\mathbf{c} - \mathbf{c}^*\|^{-1} \left[\|c_i - c_i^*\|^2 \right. \\ &\quad \left. + 2 \left\langle x - \frac{c_i + c_j}{2}, c_i - c_j \right\rangle \right], \end{aligned}$$

where we recall that $W_i(\mathbf{c})$ denotes an element of a Voronoi partition, such that $\bar{W}_i(\mathbf{c}) \subset V_i(\mathbf{c})$. The proof of (7) consists in applying Theorem 4.1 to the two terms in the right-hand side of (13).

The first term on the right-hand side of (13) may be thought of as the dominant term in the decomposition of the loss. Define

$$\mathcal{F}_2 = \{ \langle \mathbf{c} - \mathbf{c}^*(\mathbf{c}), \Delta(\mathbf{c}^*(\mathbf{c}), \cdot) \rangle \mid \mathbf{c} \in \mathcal{B}(0, M) \}.$$

In order to apply Theorem 4.1, the following lemmas are needed.

LEMMA 4.2. *Let $f \in \mathcal{F}_2$, then*

$$\begin{cases} \|f\|_\infty & \leq 8M \\ \text{Var}_P(f) & \leq 4\|\mathbf{c} - \mathbf{c}^*(\mathbf{c})\|^2 R(\mathbf{c}^*). \end{cases}$$

PROOF OF LEMMA 4.2. Elementary calculation shows that

$$\begin{aligned} \text{Var}(\langle \mathbf{c} - \mathbf{c}^*, \Delta(\mathbf{c}^*, \cdot) \rangle) &= P(\langle \mathbf{c} - \mathbf{c}^*, \Delta(\mathbf{c}^*, \cdot) \rangle)^2 - (P(\langle \mathbf{c} - \mathbf{c}^*, \Delta(\mathbf{c}^*, \cdot) \rangle))^2 \\ &= \sum_{i=1}^k P \left[\langle c_i - c_i^*, -2(x - c_i^*) \rangle^2 1_{V_i(\mathbf{c}^*)}(x) \right] \\ &\leq 4\|\mathbf{c} - \mathbf{c}^*\|^2 R(\mathbf{c}^*). \end{aligned}$$

□

Let $\omega_2(f)$ be defined as $4\|\mathbf{c} - \mathbf{c}^*(\mathbf{c})\|^2 R(\mathbf{c}^*)$. It remains to bound from above the expectation of the maximum deviation between P and P_n over the set \mathcal{F}_2 .

LEMMA 4.3. *One has*

$$(14) \quad \mathbb{E} \sup_{f \in \mathcal{F}_2, \omega_2(f) \leq \delta} |(P - P_n)f| \leq \frac{2|\mathcal{M}|}{\sqrt{n}} \sqrt{\delta}$$

PROOF. This proof is inspired from the proof of Lemma 4.3 in [5]. The first step is the following

$$\begin{aligned} \mathbb{E} \sup_{f \in \mathcal{F}_2, \omega_2(f) \leq \delta} |(P - P_n)f| \\ \leq \mathbb{E} \sup_{\mathbf{c}^* \in \mathcal{M}, \|\mathbf{c} - \mathbf{c}^*\| \leq \sqrt{\delta/4R(\mathbf{c}^*)}} |(P - P_n) \langle \mathbf{c} - \mathbf{c}^*, \Delta(\mathbf{c}^*, \cdot) \rangle|. \end{aligned}$$

For a general function $h(Z)$ depending on a random map Z , we denote by $\mathbb{E}_Z h$ the expectation of h taken with respect to Z . Introducing some

Rademacher independent random variables σ_i and using a symmetrization inequality such as in Section 2.2 of [11] leads to

$$\begin{aligned} \mathbb{E} \sup_{\|\mathbf{c}-\mathbf{c}^*\| \leq \sqrt{\delta/4R(\mathbf{c}^*)}, \mathbf{c}^* \in \mathcal{M}} |(P - P_n) \langle \mathbf{c} - \mathbf{c}^*, \Delta(\mathbf{c}^*, \cdot) \rangle| \\ \leq 2\mathbb{E}_X \mathbb{E}_\sigma \sup_{\|\mathbf{c}-\mathbf{c}^*\| \leq \sqrt{\delta/4R(\mathbf{c}^*)}, \mathbf{c}^* \in \mathcal{M}} \left\langle \mathbf{c} - \mathbf{c}^*, \frac{1}{n} \sum_{i=1}^n \sigma_i \Delta(\mathbf{c}^*, X_i) \right\rangle \\ \leq \sqrt{\delta/4R(\mathbf{c}^*)} 2\mathbb{E}_X \mathbb{E}_\sigma \sup_{\mathbf{c}^* \in \mathcal{M}} \left\| \frac{1}{n} \sum_{i=1}^n \sigma_i \Delta(\mathbf{c}^*, X_i) \right\|, \end{aligned}$$

using Cauchy-Schwarz inequality. Eventually,

$$\begin{aligned} \mathbb{E}_X \mathbb{E}_\sigma \sup_{\mathbf{c}^* \in \mathcal{M}} \left\| \frac{1}{n} \sum_{i=1}^n \sigma_i \Delta(\mathbf{c}^*, X_i) \right\| &\leq \sum_{\mathbf{c}^* \in \mathcal{M}} \mathbb{E}_X \mathbb{E}_\sigma \left\| \frac{1}{n} \sum_{i=1}^n \sigma_i \Delta(\mathbf{c}^*, X_i) \right\| \\ &\leq \sum_{\mathbf{c}^* \in \mathcal{M}} \sqrt{\mathbb{E}_X \mathbb{E}_\sigma \left\| \frac{1}{n} \sum_{i=1}^n \sigma_i \Delta(\mathbf{c}^*, X_i) \right\|^2} \\ &\leq \sum_{\mathbf{c}^* \in \mathcal{M}} \frac{1}{\sqrt{n}} \sqrt{\mathbb{E}_X \|\Delta(\mathbf{c}^*, X)\|^2} \\ &\leq \frac{2|\mathcal{M}| \sqrt{R(\mathbf{c}^*)}}{\sqrt{n}}, \end{aligned}$$

where Jensen's inequality has been used to obtain the second line. This gives the desired result. \square

The contribution of the first term in the right-hand side of (13) is described by the following proposition.

PROPOSITION 4.2. *Let K_2 be a positive constant and $x > 0$. Then, with probability larger than $1 - e^{-x}$*

$$\begin{aligned} (P - P_n) \langle \hat{\mathbf{c}}_n - \mathbf{c}^*(\hat{\mathbf{c}}_n), \Delta(\mathbf{c}^*(\hat{\mathbf{c}}_n), \cdot) \rangle &\leq K_2^{-1} \left[4R(\mathbf{c}^*) \|\hat{\mathbf{c}}_n - \mathbf{c}^*(\hat{\mathbf{c}}_n)\|^2 \right. \\ &\quad \left. + \frac{48^2 |\mathcal{M}|^2 K_2^2}{n} + \frac{9K_2^2 + 128M^2 \sqrt{k} K_2}{4n} x \right]. \end{aligned}$$

The proof follows from a direct application of Theorem 4.1 to the set \mathcal{F}_2 , replacing the value $C(1/2)$ with 4 to ease the calculation.

The second term in the right-hand side of (13) may be thought of as a residual term. Deriving sharper bounds on this term requires more accurate chaining techniques, as exposed below. Define

$$\mathcal{F}_3 = \{ \|\mathbf{c} - \mathbf{c}^*(\mathbf{c})\| R(\mathbf{c}, \mathbf{c}^*(\mathbf{c}), \cdot) \mid \mathbf{c} \in \mathcal{B}(0, M) \}.$$

In order to apply Theorem 4.1, the following intermediate results are needed.

LEMMA 4.4. *Let $f \in \mathcal{F}_3$, then*

$$\begin{cases} \|f\|_\infty & \leq 2M\sqrt{k}C_\infty \\ \text{Var}_P(f) & \leq C_\infty^2 \|\mathbf{c} - \mathbf{c}^*(\mathbf{c})\|^2, \end{cases}$$

with

$$C_\infty = (2\sqrt{k} + 8\sqrt{2})M.$$

PROOF OF LEMMA 4.4. The proof of Lemma 4.4 follows from a bound on $R(\mathbf{c}, \mathbf{c}^*(\mathbf{c}), x)$, namely

$$\begin{aligned} |R(\mathbf{c}, \mathbf{c}^*, x)| &= \|\mathbf{c} - \mathbf{c}^*\|^{-1} \left| \sum_{i,j} 1_{V_i(\mathbf{c}^*) \cap W_j(\mathbf{c})} \left(\|c_i - c_i^*\|^2 \right. \right. \\ &\quad \left. \left. + 2 \left\langle x - \frac{c_i + c_j}{2}, c_i - c_j \right\rangle \right) \right| \\ &\leq \|\mathbf{c} - \mathbf{c}^*\|^{-1} \left[\sum_i \|c_i - c_i^*\|^2 1_{V_i(\mathbf{c}^*)} \right. \\ &\quad \left. + \sum_{i \neq j} 2 \left| \left\langle x - \frac{c_i + c_j}{2}, c_i - c_j \right\rangle \right| 1_{V_i(\mathbf{c}^*) \cap W_j(\mathbf{c})} \right]. \end{aligned}$$

Since, for all j in $\{1, \dots, k\}$, $W_j(\mathbf{c}) \subset V_j(\mathbf{c})$, applying Lemma 4.1 leads to

$$\begin{aligned} |R(\mathbf{c}, \mathbf{c}^*, x)| &\leq \|\mathbf{c} - \mathbf{c}^*\|^{-1} \left[\|\mathbf{c} - \mathbf{c}^*\|^2 + 8\sqrt{2}M \|\mathbf{c} - \mathbf{c}^*\| 1_{N^*(\frac{4\sqrt{2}M}{B} \|\mathbf{c} - \mathbf{c}^*\|)} \right] \\ (15) \quad &\leq \|\mathbf{c} - \mathbf{c}^*\| + 8\sqrt{2}M 1_{N^*(\frac{4\sqrt{2}M}{B} \|\mathbf{c} - \mathbf{c}^*\|)} := F_{\|\mathbf{c} - \mathbf{c}^*\|}(x). \end{aligned}$$

Elementary calculations show that, for any $\delta > 0$,

$$\|F_\delta\|_\infty \leq (2\sqrt{k} + 8\sqrt{2})M = C_\infty,$$

from which we deduce the desired upper bounds on $\text{Var}_P(f)$ and $\|f\|_\infty$, for f in \mathcal{F}_3 . \square

Let $\omega_3(f)$ be defined as $C_\infty^2 \|\mathbf{c} - \mathbf{c}^*(\mathbf{c})\|^2$. The complexity term associated with the class of functions \mathcal{F}_3 can be bounded as follows.

PROPOSITION 4.3.

$$\mathbb{E} \sup_{f \in \mathcal{F}_3, \omega_3(f) \leq \delta} |(P - P_n)f| \leq \frac{8Q(k, d)}{\sqrt{C_\infty n}} \sqrt{C_2(\sqrt{\delta}/C_\infty)\sqrt{\delta}},$$

where

$$\begin{cases} C_2(r) &= r + 8\sqrt{2}Mp \left(\frac{4\sqrt{2}M}{B}r \right) \\ Q(k, d) &= 8\sqrt{K_0 P(k, d) \log(k^2(4k-2))} \\ P(k, d) &= k^2(2(k-1)(d+1) + 24(3d+4)), \end{cases}$$

and K_0 is an absolute constant.

The proof of Proposition 4.3 is based on a result of Mendelson and Vershynin in [18] and its application to a more accurate version of Dudley's integral. For clarity, the proof is postponed to Section 5.2. Since P satisfies a margin condition with parameters (r_0, κ) , with $\kappa \leq \frac{Bp_{\min}}{128M^2}$, considering the two cases $4\sqrt{2}Mr/B \leq r_0$ and $4\sqrt{2}Mr/B \geq r_0$ yields that

$$8\sqrt{2}Mp \left(\frac{4\sqrt{2}M}{B}r \right) \leq \frac{64M^2}{Br_0}r,$$

for $r \geq 0$. Using this inequality to bound C_2 from above in Proposition 4.3 leads to the following complexity result

$$(16) \quad \mathbb{E} \sup_{f \in \mathcal{F}_3, \omega_3(f) \leq \delta} |(P - P_n)f| \leq \frac{\Xi_3}{\sqrt{n}} \delta^{\frac{3}{4}},$$

where

$$\Xi_3 = \frac{8MQ(k, d)}{C_\infty \sqrt{Br_0}}.$$

Let Φ_3 be defined as $\frac{\Xi_3}{\sqrt{n}} \delta^{\frac{3}{4}}$. Remark that Φ_3 is a sub-3/4 function. Consequently, for any $D > 0$, the solution of the equation $\Phi_3(\delta) = \delta/D$ is

$$\delta_3^* = \frac{(D\Xi_3)^4}{n^2}.$$

Choosing $K_3 > 0$, and $D = 6K_3C(3/4)$ in Theorem 4.1 and taking into account that $C(3/4) \leq 10$ leads to the following proposition.

PROPOSITION 4.4. *Let $K_3 > 0$. Then, with probability larger than $1 - e^{-x}$,*

$$(P - P_n)(\|\hat{\mathbf{c}}_n - \mathbf{c}^*(\hat{\mathbf{c}}_n)\|R(\hat{\mathbf{c}}_n, \mathbf{c}^*(\hat{\mathbf{c}}_n), \cdot)) \leq K_3^{-1} \left[C_\infty^2 \|\hat{\mathbf{c}}_n - \mathbf{c}^*(\hat{\mathbf{c}}_n)\|^2 + \frac{60^4 K_3^4 \Xi_3^4}{n^2} + \frac{9K_3^2 + 32M\sqrt{k}C_\infty K_3}{4n} x \right],$$

with $\Xi_3 = \frac{8MQ(k,d)}{C_\infty \sqrt{Br_0}}$, and Q is a function composed of products of square roots of polynomial functions in k , d , and $\log(k)$.

We are now in position to prove (7). Proposition 2.1 provides κ_0 such that

$$\kappa_0 \ell(\mathbf{c}, \mathbf{c}^*) \geq \|\mathbf{c} - \mathbf{c}^*(\mathbf{c})\|^2.$$

Choosing $K_2 = 8R(\mathbf{c}^*)\kappa_0$ in Proposition 4.2, $K_3 = 2C_\infty^2 \kappa_0$ in Proposition 4.4 and summing the two resulting inequalities leads to (7), valid on a set which has probability larger than $1 - 2e^{-x}$.

4.3. *Proof of Proposition 3.1.* Throughout this subsection, for a codebook \mathbf{c} , let Q denote the associated nearest-neighbor quantizer. In the general case, such an association depends on how the boundaries are allocated. However, since the distributions involved in the minimax result have a density, how boundaries are allocated will not matter.

Let $k \geq 3$ be an integer. For convenience k is assumed to be divisible by 3. Let $m = 2k/3$. Let z_1, \dots, z_m denote a 6Δ -net in $\mathcal{B}(0, M - \rho)$, where $\Delta > 0$, and w_1, \dots, w_m a sequence of vectors such that $\|w_i\| = \Delta$. Finally, denote by U_i the ball $\mathcal{B}(z_i, \rho)$ and by U'_i the ball $\mathcal{B}(z_i, \rho)$. Slightly anticipating, define $\rho = \frac{\Delta}{16}$.

To get the largest Δ such that for all $i = 1, \dots, k$ U_i and U'_i are included in $\mathcal{B}(0, M)$, it suffices to get the largest Δ such that there exists a 6Δ -net in $\mathcal{B}(0, M - \Delta/16)$. Since the cardinal of a 6Δ -net is larger than the largest number of balls of radius 6Δ which can be packed into $\mathcal{B}(0, M - \Delta/16)$, a sufficient condition on Δ to guarantee that a 6Δ -net can be found is given by

$$m \leq \left(\frac{M - \Delta/16}{6\Delta} \right)^d.$$

Since $\Delta \leq M$, Δ can be chosen as

$$\Delta = \frac{15M}{96m^{1/d}}.$$

For such a Δ , ρ takes the value $\rho = \frac{\Delta}{16} = \frac{15M}{1536m^{1/d}}$. Therefore, it only depends on k , d , and M .

Let $z = (z_i)_{i=1,\dots,m}$ and $w = (w_i)_{i=1,\dots,m}$ be sequences as described above, such that, for $i = 1, \dots, k$, U_i and U'_i are included in $\mathcal{B}(0, M)$. For a fixed $\sigma \in \{-1, +1\}^m$ such that $\sum_{i=1}^m \sigma_i = 0$, let P_σ be defined as

$$\begin{cases} P_\sigma(U_i) &= \frac{1+\sigma_i\delta}{2m} \\ P_\sigma(U'_i) &= \frac{1+\sigma_i\delta}{2m} \\ P_\sigma &\underset{U_i}{\sim} (\rho - \|x - z_i\|)1_{\|x - z_i\| \leq \rho} d\lambda(x) \\ P_\sigma &\underset{U'_i}{\sim} (\rho - \|x - z_i - w_i\|)1_{\|x - z_i - w_i\| \leq \rho} d\lambda(x), \end{cases}$$

where λ denote the Lebesgue measure. These cone-shaped distributions are designed to have a continuous density, as done in [1]. To be more precise, for τ in $\{-1, +1\}^{\frac{m}{2}}$, $\sigma(\tau)$ is defined as the sequence in $\{-1, +1\}^m$ such that

$$\begin{cases} \sigma_i(\tau) &= \tau_i \\ \sigma_{i+\frac{m}{2}}(\tau) &= -\sigma_i(\tau), \end{cases}$$

for $i = 1, \dots, \frac{m}{2}$. Finally, for a quantizer Q let $R(Q, P_\sigma)$ denote the distortion of Q in the case where the source distribution is P_σ .

Similarly, for σ in $\{-1, +1\}^m$ satisfying $\sum_{i=1}^m \sigma_i = 0$, let Q_σ denote the quantizer defined by $Q_\sigma(U_i) = Q_\sigma(U'_i) = z_i + \omega_i/2$ if $\sigma_i = -1$, $Q_\sigma(U_i) = z_i$ and $Q_\sigma(U'_i) = z_i + \omega_i$ if $\sigma_i = +1$. Let \mathcal{Q} denote the set of such quantizers. It can be proved that only quantizers in \mathcal{Q} have to be considered.

PROPOSITION 4.5. *Assume that $A \geq 6$, $\delta \leq 1/3$, $\Delta > 0$, and $\rho \leq \frac{\Delta}{16}$. Then, for every quantizer Q there exists a quantizer Q_σ in \mathcal{Q} such that*

$$\forall P_{\sigma'} \quad R(Q_\sigma, P_{\sigma'}) \leq R(Q, P_{\sigma'}).$$

The proof of Proposition 4.5 follows the proof of Step 3 of Theorem 1 in [3], replacing distributions supported on a finite set with distributions supported on small balls. Provided that the radius of these balls are small enough, the results are nearly the same in the two cases. The proof of Proposition 4.5 is given in Section 5.4.

For any σ and σ' in $\{-1, 1\}^m$, denote by $\rho(\sigma, \sigma') = \sum_{i=1}^m |\sigma_i - \sigma'_i|$, and by $H(P_\sigma, P_{\sigma'})$ the Hellinger distance between P_σ and $P_{\sigma'}$. To apply Assouad's Lemma to the set $\{P_{\sigma(\tau)}\}_{\tau \in \{-1, +1\}^{\frac{m}{2}}}$, the following lemma is needed:

LEMMA 4.5. *Let τ and τ' denote two sequences in $\{-1, +1\}^{\frac{m}{2}}$ such that $\rho(\tau, \tau') = 2$, then*

$$H(P_{\sigma(\tau)}^{\otimes n}, P_{\sigma(\tau')}^{\otimes n}) \leq \frac{4n\delta^2}{m},$$

where $P^{\otimes n}$ denotes the product law of a n -sample drawn from P .

Furthermore, for any σ and σ' in $\{-1, +1\}^m$,

$$R(Q_{\sigma'}, P_{\sigma}) = R(Q_{\sigma}, P_{\sigma}) + \frac{\Delta^2 \delta}{8m} \rho(\sigma, \sigma').$$

Equipped with Lemma 4.5, a direct application of Assouad's Lemma as in Theorem 2.12 of [21] yields that, provided that $\delta = \frac{\sqrt{m}}{2\sqrt{n}}$,

$$\mathbb{E} \sup_{\tau \in \{-1, +1\}^{\frac{m}{2}}} R(\hat{Q}_n, P_{\sigma(\tau)}) - R(Q_{\sigma(\tau)}, P_{\sigma(\tau)}) \geq c_0 M^2 \sqrt{\frac{k^{1-\frac{4}{d}}}{n}},$$

for any empirically designed quantizer \hat{Q}_n , where c_0 is an explicit constant.

Finally, it may be noticed that, for every $\delta \leq \frac{1}{3}$ and σ , P_{σ} satisfies a margin condition as in (8), and is ε -separated, with

$$\varepsilon = \frac{\Delta^2 \delta}{2m}.$$

This concludes the proof of Proposition 3.1.

As a remark, it is worth mentioning that, whenever $\varepsilon \sim 1/n^\alpha$, with $\alpha < 1/2$, no interesting minimax lower bound can be derived using the distributions $\{P_{\sigma}\}$'s. In fact, it can be proved, making use of the approach exposed in [3], that the empirical risk minimization strategy \hat{Q}_n achieves the uniform rate

$$\sup_{\{P_{\sigma}\}_{\sigma \in \{-1, +1\}^m}} R(\hat{Q}_n, P_{\sigma}) - R(Q_{\sigma}, P_{\sigma}) \leq C(M, k, d) n^{-\alpha} e^{-c(k)n^{1-2\alpha}},$$

where $C(M, k, d)$ and $c(k)$ are constants. Consequently, in order to get a minimax lower bound matching the bound offered in Theorem 3.1, more general probability distributions should be used.

4.4. *Proof of Proposition 3.2.* As mentioned below Proposition 3.2, the inequality

$$\frac{\theta_{\min}}{\theta_{\max}} \geq \frac{2048k\sigma^2}{(1-\varepsilon)\tilde{B}^2(1-e^{-\tilde{B}^2/2048\sigma^2})},$$

ensures that, for every j in $\{1, \dots, k\}$, there exists i in $\{1, \dots, k\}$ such that $\|c_i^* - m_j\| \leq \tilde{B}/16$. To be more precise, let \mathbf{m} denote the vector of means (m_1, \dots, m_k) , then

$$\begin{aligned} R(\mathbf{m}) &\leq \sum_{i=1}^k \frac{\theta_i}{2\pi\sigma^2 N_i} \int_{V_i(\mathbf{m})} \|x - m_i\|^2 e^{-\frac{\|x - m_i\|^2}{2\sigma^2}} dx \\ &\leq \frac{p_{max}}{2(1-\varepsilon)\pi\sigma^2} \sum_{i=1}^k \int_{\mathbb{R}^2} \|x - m_i\|^2 e^{-\frac{\|x - m_i\|^2}{2\sigma^2}} dx \\ &\leq \frac{2kp_{max}\sigma^2}{1-\varepsilon}. \end{aligned}$$

Assume that there exists i in $\{1, \dots, k\}$ such that, for all j , $\|c_j^* - m_i\| \geq \tilde{B}/16$. Then

$$\begin{aligned} R(\mathbf{c}) &\geq \frac{\theta_i}{2\pi\sigma^2} \int_{B(m_i, \tilde{B}/32)} \frac{\tilde{B}^2}{1024} e^{-\frac{\|x - m_i\|^2}{2\sigma^2}} \\ &\geq \frac{\tilde{B}^2 \theta_{min}}{2048\pi\sigma^2} \int_{B(m_i, \tilde{B}/32)} e^{-\frac{\|x - m_i\|^2}{2\sigma^2}} \\ &> \frac{\tilde{B}^2 \theta_{min}}{1024} \left(1 - e^{-\frac{\tilde{B}^2}{2048\sigma^2}}\right) \\ &> R(\mathbf{m}). \end{aligned}$$

Hence the contradiction. Up to relabeling, it is now assumed that for $i = 1, \dots, k$, $\|m_i - c_i^*\| \leq \tilde{B}/16$. Take y in $N^*(x)$, for $x \leq \frac{\tilde{B}}{8}$, then, for every i in $\{1, \dots, k\}$,

$$\|y - m_i\| \geq \frac{\tilde{B}}{4},$$

which leads to

$$\sum_{i=1}^k \frac{\theta_i}{2\pi\sigma^2 N_i} \|y - m_i\|^2 e^{-\frac{\|y - m_i\|^2}{2\sigma^2}} \leq \frac{k\theta_{max}}{(1-\varepsilon)2\pi\sigma^2} e^{-\frac{\tilde{B}^2}{32\sigma^2}}.$$

Since the Lebesgue measure of $N^*(x)$ is smaller than $4k\pi Mx$, it follows that

$$P(N^*(x)) \leq \frac{2k^2 M \theta_{max}}{(1-\varepsilon)\sigma^2} e^{-\frac{\tilde{B}^2}{32\sigma^2}} x.$$

On the other hand, $\|m_i - c_i^*\| \leq \tilde{B}/16$ yields that

$$B(m_i, 3\tilde{B}/8) \subset V_i(\mathbf{c}^*).$$

Therefore,

$$\begin{aligned} P(V_i(\mathbf{c}^*)) &\geq \frac{\theta_i}{2\pi\sigma^2 N_i} \int_{\mathcal{B}(m_i, 3\tilde{B}/8)} e^{-\frac{\|x-m_i\|^2}{2\sigma^2}} dx \\ &\geq \theta_i \left(1 - e^{-\frac{9\tilde{B}^2}{128\sigma^2}}\right), \end{aligned}$$

hence $p_{\min} \geq \theta_{\min} \left(1 - e^{-\frac{9\tilde{B}^2}{128\sigma^2}}\right)$. Consequently, provided that

$$\frac{\theta_{\min}}{\theta_{\max}} \geq \frac{2048k^2 M^3}{(1-\varepsilon)7\sigma^2 \tilde{B}(e^{\tilde{B}^2/32\sigma^2} - 1)},$$

direct calculation shows that

$$P(N^*(x)) \leq \frac{Bp_{\min}}{128M^2}x.$$

5. Technical results.

5.1. *Proof of Proposition 4.1.* The proof of Proposition 4.1 is derived from the proof of Lemma 3 in [7]. Let \mathbf{c}^* be an optimal codebook, and \mathbf{c} be a codebook. We denote by $f_{\mathbf{c}^*, \mathbf{c}}$ the function $\gamma(\mathbf{c}, \cdot) - \gamma(\mathbf{c}^*, \cdot)$, so that

$$\mathcal{F}_1 = \left\{ f_{\mathbf{c}^*(\mathbf{c}), \mathbf{c}} \mid \mathbf{c} \in \mathcal{B}(0, M)^k \right\}.$$

Let $\Psi_1(r)$ denote the function

$$\Psi_1(r) = \mathbb{E} \sup_{\mathbf{c}^* \in \mathcal{M}, \|\mathbf{c} - \mathbf{c}^*(\mathbf{c})\| \leq r} |(P - P_n)f_{\mathbf{c}^*(\mathbf{c}), \mathbf{c}}|.$$

Since

$$\{(\mathbf{c}^*(\mathbf{c}), \mathbf{c}^*) \mid \|\mathbf{c} - \mathbf{c}^*(\mathbf{c})\| \leq r\} \subset \{(\mathbf{c}^*, \mathbf{c}) \mid \mathbf{c}^* \in \mathcal{M}, \|\mathbf{c} - \mathbf{c}^*\| \leq r\},$$

it is easy to see that

$$\Psi_1(r) \leq \mathbb{E} \sup_{\mathbf{c}^* \in \mathcal{M}, \|\mathbf{c} - \mathbf{c}^*\| \leq r} |(P - P_n)f_{\mathbf{c}^*, \mathbf{c}}|.$$

The rest of the proof is derived from a chaining technique, used in the proof of Proposition 5.1 in [12] or Lemma 3 in [7]. Set $\varepsilon_j = 2^{-j}r$, for $j \geq 0$, and for every \mathbf{c}^* in \mathcal{M} , denote by $N_j(\mathbf{c}^*)$ an ε_j net of $\mathcal{B}(\mathbf{c}^*, r)$, such that for every \mathbf{c} in $\mathcal{B}(\mathbf{c}^*, r)$ there exists \mathbf{c}_j in $N_j(\mathbf{c}^*)$ such that $\|\mathbf{c}_j - \mathbf{c}\| \leq \varepsilon_j$.

According to the proof of Theorem 2 in [2] or Lemma 3 in [7], such an $N_j(\mathbf{c}^*)$ can be defined, with

$$|N_j(\mathbf{c}^*)| \leq \left(\frac{2r\sqrt{kd}}{\varepsilon_j} \right) := n(\varepsilon_j).$$

By a dominated convergence Theorem, for any fixed \mathbf{c}^* in \mathcal{M} and \mathbf{c} ,

$$f_{\mathbf{c}^*, \mathbf{c}_j} \xrightarrow[j \rightarrow \infty]{L_1(P), a.s.} f_{\mathbf{c}^*, \mathbf{c}}.$$

This allows us to decompose the expression of Ψ_1 as follows.

$$\begin{aligned} \Psi_1(r) &\leq \mathbb{E} \sup_{\mathbf{c}^* \in \mathcal{M}, \|\mathbf{c} - \mathbf{c}^*\| \leq r} |(P - P_n)f_{\mathbf{c}^*, \mathbf{c}}| \\ &\leq \mathbb{E} \sup_{\mathbf{c}^* \in \mathcal{M}, \mathbf{c}_0 \in N_0(\mathbf{c}^*)} |(P - P_n)f_{\mathbf{c}^*, \mathbf{c}_0}| \\ &\quad + \sum_{j>1} \mathbb{E} \sup_{\mathbf{c}^* \in \mathcal{M}, \mathbf{c}_j \in N_j(\mathbf{c}^*), \mathbf{c}_{j-1} \in N_{j-1}(\mathbf{c}^*)} |(P - P_n)(f_{\mathbf{c}^*, \mathbf{c}_j} - f_{\mathbf{c}^*, \mathbf{c}_{j-1}})|, \\ &:= A_1 + A_2. \end{aligned}$$

It remains to bound from above these two terms.

Bound on A_1

Introducing some Rademacher random variables $\sigma_i, i = 1, \dots, n$ and using the symmetrization principle as in [11] leads to

$$\begin{aligned} \mathbb{E} \sup_{\mathbf{c}^* \in \mathcal{M}, \mathbf{c}_0 \in N_0(\mathbf{c}^*)} |(P - P_n)f_{\mathbf{c}^*, \mathbf{c}_0}| \\ \leq 2\mathbb{E}_X \mathbb{E}_\sigma \sup_{\lambda = \pm 1, \mathbf{c}^* \in \mathcal{M}, \mathbf{c}_0 \in N_0(\mathbf{c}^*)} \frac{1}{n} \sum_{i=1}^n \sigma_i \lambda f_{\mathbf{c}^*, \mathbf{c}_0}(X_i). \end{aligned}$$

Let introduce here a maximal inequality derived from Lemma 2.3 in [16].

LEMMA 5.1. *Let x_1, \dots, x_n denote a sequence of points in \mathcal{X} , and let $\sigma_1, \dots, \sigma_n$ denote a sequence of independent Rademacher random variables. Let \mathcal{F} be a set of real valued functions over \mathcal{X} such that $|\mathcal{F}| < \infty$, and*

$$\sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n f^2(x_i) \leq v.$$

Then

$$\mathbb{E}_\sigma \sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i f(x_i) \leq \sqrt{2v \log(|\mathcal{F}|)}.$$

In our case, for all \mathbf{c}^* in \mathcal{M} and \mathbf{c}_0 in $N_0(\mathbf{c}^*)$,

$$\frac{1}{n} \sum_{i=1}^n f_{\mathbf{c}^*, \mathbf{c}_0}^2(X_i) \leq \frac{16M^2 r^2}{n},$$

and

$$|\{\lambda = \pm 1, \mathbf{c}^* \in \mathcal{M}, \mathbf{c}_0 \in N_0(\mathbf{c}^*)\}| \leq |\mathcal{M}|(4\sqrt{kd})^{kd}.$$

Therefore, a direct application of Lemma 5.1 yields that

$$A_1 \leq \frac{8\sqrt{2}M}{\sqrt{n}} \sqrt{kd \log(4|\mathcal{M}|\sqrt{kd})}.$$

Bound on A_2

Let $j > 1$. Using the same symmetrization argument as above leads to

$$\begin{aligned} A_{2,j} &:= \mathbb{E} \sup_{\mathbf{c}^* \in \mathcal{M}, \mathbf{c}_j \in N_j(\mathbf{c}^*), \mathbf{c}_{j-1} \in N_{j-1}(\mathbf{c}^*)} |(P - P_n)(f_{\mathbf{c}^*, \mathbf{c}_j} - f_{\mathbf{c}^*, \mathbf{c}_{j-1}})| \\ &\leq \mathbb{E}_X \mathbb{E}_\sigma \sup_{\substack{\lambda = \pm 1, \mathbf{c}^* \in \mathcal{M}, \\ \mathbf{c}_j \in N_j(\mathbf{c}^*), \mathbf{c}_{j-1} \in N_{j-1}(\mathbf{c}^*)}} \frac{1}{n} \sum_{i=1}^n \sigma_i \lambda (f_{\mathbf{c}^*, \mathbf{c}_j}(X_i) - f_{\mathbf{c}^*, \mathbf{c}_{j-1}}(X_i)). \end{aligned}$$

Since

$$\|f_{\mathbf{c}^*, \mathbf{c}_j}(\cdot) - f_{\mathbf{c}^*, \mathbf{c}_{j-1}}(\cdot)\|_\infty \leq 8Mr2^{-(j-1)},$$

and

$$|\{\lambda = \pm 1, \mathbf{c}^* \in \mathcal{M}, \mathbf{c}_j \in N_j(\mathbf{c}^*), \mathbf{c}_{j-1} \in N_{j-1}(\mathbf{c}^*)\}| \leq 2|\mathcal{M}|n(\varepsilon_j)^2,$$

a direct application of Lemma 5.1 leads to

$$A_{2,j} \leq 64Mr \sqrt{kd \log(|\mathcal{M}|\sqrt{kd}2^{j+2})} 2^{-(j-1)}.$$

Comparing a sum with an integral, and observing that

$$\int_0^1 \sqrt{\log(-x)} dx \leq 1$$

ensures that

$$A_2 = \sum_{j>1} A_{2,j} \leq \frac{256Mr}{\sqrt{n}} \left(\sqrt{\log(|\mathcal{M}|\sqrt{kd})} + 1 \right).$$

Combining the two bounds and remarking that

$$\mathbb{E} \sup_{f \in \mathcal{F}_1, \omega_1(f) \leq \delta} |(P - P_n)f| \leq \Psi_1 \left(\frac{\sqrt{\delta}}{4M} \right)$$

gives the result of Proposition 4.1.

5.2. *Proof of Proposition 4.3.* The proof of Proposition 4.3 is based on a sharper chaining technique than the one used in Proposition 5.1 in [12]. We intend to bound from above the complexity term

$$\mathbb{E} \sup_{\omega_3(f) \leq \delta, f \in \mathcal{F}_3} |(P - P_n)f|.$$

To this aim, define

$$\Psi_3(r) = \mathbb{E} \sup_{\|\mathbf{c} - \mathbf{c}^*(\mathbf{c})\| \leq r} |(P - P_n)R(\mathbf{c}, \mathbf{c}^*, \cdot)|,$$

where we recall that

$$R(\mathbf{c}, \mathbf{c}^*, x) = \sum_{i,j=1,\dots,k} 1_{V_i(\mathbf{c}^*) \cap W_j(\mathbf{c})} \|\mathbf{c} - \mathbf{c}^*\|^{-1} \left[\|c_i - c_i^*\|^2 + 2 \left\langle x - \frac{c_i + c_j}{2}, c_i - c_j \right\rangle \right],$$

where $(W_1(\mathbf{c}), \dots, W_k(\mathbf{c}))$ is a Voronoi partition, defined in Section 2. For technical reasons, this Voronoi partition must be specified. Denote by $\mathcal{C}(p)$ the set of subsets of \mathbb{R}^d made of intersections of at most p half spaces (closed or open).

Since $R(\mathbf{c}, \mathbf{c}^*, \cdot)$ does not depend on how ties are broken, or, in other words, $R(\mathbf{c}, \mathbf{c}^*, \cdot)$ does not depend on the choice of $W_j(\mathbf{c})$'s among the partition cells satisfying

$$\overset{o}{V}_j(\mathbf{c}) \subset W_j(\mathbf{c}) \subset V_j(\mathbf{c}),$$

we choose a Voronoi partition such that every $W_j \in \mathcal{C}(k-1)$. For instance, if $H_{i,j}$ denotes the closed half-space $\{\|x - c_i\| \leq \|x - c_j\|\}$ and $\overset{o}{H}_{i,j}$ the open half space $\{\|x - c_i\| < \|x - c_j\|\}$. It is possible to build a Voronoi partition such that every cell is in $\mathcal{C}(k-1)$, choosing

$$W_j(\mathbf{c}) = \bigcap_{i < j} \overset{o}{H}_{i,j} \cap \bigcap_{i > j} H_{i,j}.$$

In short, this convention consists in allocating points on boundaries between V_j 's to the smallest possible index. As a consequence, it is immediate that

$$\Psi_3(r) = \mathbb{E} \sup_{\|\mathbf{c} - \mathbf{c}^*(\mathbf{c})\| \leq r, W_j(\mathbf{c}) \in \mathcal{C}(k-1)} |(P - P_n)R(\mathbf{c}, \mathbf{c}^*, \cdot)|.$$

The following set of function of interest is then introduced.

$$\mathcal{G}(r) = \{0\} \cup \frac{1}{F_r} \left\{ \lambda \|\mathbf{c} - \mathbf{c}^*\|^{-1} \sum_{i,j} 1_{V_i(\mathbf{c}^*) \cap W_j(\mathbf{c}) \cap \mathcal{B}(0,M)} \left(\|c_i - c_i^*\|^2 \right. \right. \\ \left. \left. + 2 \left\langle x - \frac{c_i + c_j}{2}, c_i - c_j \right\rangle \right) \mid \lambda = \pm 1, \mathbf{c} \in \mathcal{B}(0, M)^k, \mathbf{c}^* \in \mathcal{M}, W_j \in \mathcal{C}(k-1) \right\},$$

where F_r is defined in (15) as an envelope of $\mathcal{G}(r)$.

Let $\sigma_1, \dots, \sigma_n$ denote a sequence of independent Rademacher variables. As developed in the proof of Proposition 4.1, the first step is a symmetrization inequality

$$\begin{aligned} \Psi_3(r) &\leq \mathbb{E} \sup_{\mathbf{c}^* \in \mathcal{M}, \|\mathbf{c} - \mathbf{c}^*\| \leq r} |(P - P_n)R(\mathbf{c}, \mathbf{c}^*, \cdot)| \\ &\leq 2\mathbb{E}_X \mathbb{E}_\sigma \sup_{\lambda = \pm 1, \mathbf{c}^* \in \mathcal{M}, \|\mathbf{c} - \mathbf{c}^*\| \leq r} \frac{1}{n} \sum_{i=1}^n \lambda \sigma_i R(\mathbf{c}, \mathbf{c}^*, X_i) \\ &\leq 2\mathbb{E}_X \mathbb{E}_\sigma \sup_{g \in \mathcal{G}(r)} \frac{1}{n} \sum_{i=1}^n \sigma_i F_r(X_i) g(X_i) \\ &:= 2\mathbb{E}_X \mathcal{R}_n. \end{aligned}$$

The next step is to chain the set $\mathcal{G}(r)$. To this aim, define for any set of real valued function \mathcal{F} , any norm $\|\cdot\|$ on \mathcal{F} and any $\varepsilon > 0$, the covering number $\mathcal{N}(\mathcal{F}, \|\cdot\|, \varepsilon)$ as the cardinal of the smallest covering of \mathcal{F} with balls of radius ε for the norm $\|\cdot\|$.

To be more precise, for any g in $\mathcal{G}(r)$, and any finite subset $S \subset \mathbb{R}$, we define

$$\|g\|_{L_2(S)} = \sqrt{\frac{1}{|S|} \sum_{s \in S} g^2(s)},$$

and, with a slight abuse of notation, $\|g\|_{L_2(P_n)} = \sqrt{1/n \sum_{i=1}^n g^2(X_i)}$. The technical result concerning the covering numbers of $\mathcal{G}(r)$ is the following.

PROPOSITION 5.1. *Let S be a finite set, and $0 < \varepsilon < 1$. There exists some constant $K > 0$, not depending on S , such that*

$$\mathcal{N}(\mathcal{G}(r), \varepsilon, L_2(S)) \leq \left(\frac{k^2(4k-2)}{\varepsilon} \right)^{KP(k,d)},$$

with $P(k, d) = k^2(2(k-1)(d+1) + 24(3d+4))$.

For clarity, the proof of Proposition 5.1 is postponed to the following subsection. An immediate consequence of Proposition 5.1 is that

$$\mathcal{N}(\mathcal{G}(r), \varepsilon, L_2(P_n)) \leq \left(\frac{k^2(4k-2)}{\varepsilon} \right)^{KP(k,d)} := n(\varepsilon),$$

for any n -sample X_1, \dots, X_n . Consequently, let X_1, \dots, X_n be fixed, and set $\varepsilon_0 = 1$, $\varepsilon_j = 2^{-2j}\varepsilon_0$, for $j > 1$.

For $j = 0$, since F_r is an envelope of $\mathcal{G}(r)$, a 1 covering of $\mathcal{G}(r)$ for the $L_2(P_n)$ norm is the ball of center $g_0 = 0$ and radius 1.

For $j > 1$, Proposition 5.1 provides a ε_j covering $\mathcal{G}_j(r)$ of $\mathcal{G}(r)$ for the $L_2(P_n)$ norm with cardinality at most $n(\varepsilon_j)$. For any g in $\mathcal{G}(r)$, denote by g_j the projection of g onto this covering, so that $\|g - g_j\|_{L_2(P_n)} \leq \varepsilon_j$. For short we will write $n_j = n(\varepsilon_j)$.

It is easy to see that for every i in $\{1, \dots, n\}$,

$$g_j(X_i) \xrightarrow{j \rightarrow \infty} g(X_i).$$

Then \mathcal{R}_n may be decomposed as follows:

$$\begin{aligned} \mathcal{R}_n &= \mathbb{E}_\sigma \sup_{\lambda=\pm 1, \mathbf{c}^* \in \mathcal{M}, \|\mathbf{c}-\mathbf{c}^*\| \leq r} \frac{1}{n} \sum_{i=1}^n \lambda \sigma_i R(\mathbf{c}, \mathbf{c}^*, X_i) \\ &\leq \sum_{j>1} \mathbb{E}_\sigma \sup_{g \in \mathcal{G}} \frac{1}{n} \sum_{i=1}^n \sigma_i F_r(X_i) (g_j(X_i) - g_{j-1}(X_i)) \\ &:= \sum_{j>1} b_j. \end{aligned}$$

A direct application of Lemma 5.1 for every b_j yields that

$$\begin{aligned} b_j &\leq \frac{1}{\sqrt{n}} \sqrt{2 \sup_{g \in \mathcal{G}(r)} \|(g_j - g_{j-1})F_r\|_{L_2(P_n)}^2 \log(n_j n_{j-1})} \\ &\leq \frac{1}{\sqrt{n}} \sqrt{2 \log(n_j n_{j-1}) \sup_{g \in \mathcal{G}(r)} 2C_\infty \|F_r\|_{L_2(P_n)} \|g_j - g_{j-1}\|_{L_2(P_n)}} \\ &\leq \frac{4\sqrt{C_\infty}}{\sqrt{n}} \sqrt{\|F_r\|_{L_2(P_n)}} \sqrt{\log(n(\varepsilon_j))} \sqrt{\varepsilon_{j-1}}. \end{aligned}$$

Denote by ε'_j the quantity $\sqrt{\varepsilon_j} = 2^{-j}$. Since $x \mapsto \sqrt{\log(n(x^2))}$ is non-increasing, it is quite easy to see that

$$\frac{\sqrt{\log(n(\varepsilon_j'^2))}\varepsilon'_{j-1}}{4} = \sqrt{\log(n(\varepsilon_j'^2))}\varepsilon'_{j+1} \leq \int_{\varepsilon'_{j+1}}^{\varepsilon'_j} \sqrt{\log(n(x^2))} dx.$$

From this we deduce that

$$\begin{aligned}
\sum_{j>1} \sqrt{\varepsilon_{j-1} \log(n(\varepsilon_j))} &\leq 4 \int_0^{1/2} \sqrt{\log(n(x^2))} dx \\
&\leq \int_0^{1/2} \sqrt{KP(k, d) \log\left(\frac{k^2(4k-2)}{x^2}\right)} dx \\
&\leq 2\sqrt{KP(k, d) \log(k^2(4k-2))} \\
&\quad + 4\sqrt{KP(k, d)} \int_0^{1/2} \sqrt{\log(1/x^2)} dx.
\end{aligned}$$

Since $\int_0^{1/2} \sqrt{\log(1/x^2)} dx \leq 1$, we get

$$\sum_{j>1} \sqrt{\varepsilon_{j-1} \log(n(\varepsilon_j))} \leq 8\sqrt{KP(k, d) \log(k^2(4k-2))} := Q(k, d).$$

Thus

$$\mathcal{R}_n \leq \frac{4Q(k, d)\sqrt{C_\infty}}{\sqrt{n}} \sqrt{\|F_r\|_{L_2(P_n)}}.$$

It remains now to take expectations with respect to the n -sample X_1, \dots, X_n . Since $x \mapsto \sqrt{x}$ is a concave map,

$$\mathbb{E}_X(\sqrt{\|F_r\|_{L_2(P_n)}}) \leq \sqrt{\|F_r\|_{L_2(P)}} = \sqrt{C_2(r)}.$$

Gathering all terms leads to

$$\Psi_3(r) \leq \frac{8Q(k, d)\sqrt{C_\infty C_2(r)}}{\sqrt{n}}.$$

Substituting $\sqrt{\delta}/C_\infty$ with r gives the result of Proposition 4.3.

5.3. Proof of Proposition 5.1. Let S be a finite subset of \mathbb{R}^d , and denote by $\mathcal{C}(p)$ the set of subsets of \mathbb{R}^d which are intersections of at most p half spaces (closed or open). For short, $\mathcal{N}(\mathcal{F}, \varepsilon)$ will denote $\mathcal{N}(\mathcal{F}, L_2(S), \varepsilon)$.

The proof of Proposition 5.1 is based on the following result, Theorem 1 in [18].

THEOREM 5.1. *Let P denote a measure on Ω . Let \mathcal{F} be a set of maps from Ω into $[-1, 1]$. Then, for every $0 < t < 1$,*

$$\mathcal{N}(\mathcal{F}, t, L_2(P)) \leq \left(\frac{2}{t}\right)^{Kvc(\mathcal{F}, ct)},$$

where K and c are constants, and $vc(\mathcal{F}, ct)$ denotes the t -shattering dimension of \mathcal{F} , as defined in [18].

Remark that, for every $t > 0$, $vc(\mathcal{F}, ct) \leq d_p(\mathcal{F})$, where $d_p(\mathcal{F})$ denotes the pseudo-dimension of \mathcal{F} , that is the largest integer p such that there exists $x_1, \dots, x_p \in \Omega$, and t_1, \dots, t_p real numbers, satisfying the following property: for every $\sigma \in \{-1, 1\}^p$ there exists $f_\sigma \in \mathcal{F}$ such that, for $i = 1, \dots, p$, $\sigma_i(f_\sigma(x_i) - t_i) > 0$. As a consequence, the quantity of interest is $d_p(\mathcal{G}(r))$.

Recalling that every g in \mathcal{G} can be written

$$R(\mathbf{c}, \mathbf{c}^*, x) = \|\mathbf{c} - \mathbf{c}^*\|^{-1} \sum_{i,j} 1_{V_i(\mathbf{c}^*) \cap W_j(\mathbf{c}) \cap \mathcal{B}(0,M)} \left(\|c_i - c_i^*\|^2 + 2 \left\langle x - \frac{c_i + c_j}{2}, c_i - c_j \right\rangle \right),$$

where, for every j in $\{1, \dots, k\}$, $W_j(\mathbf{c})$ is in $\mathcal{C}(k-1)$, it may be noticed that g takes the form of a sum of k^2 maps of the type $\ell 1_C 1_{\mathcal{B}(0,1)}$, where ℓ denotes an affine map, and C is an element of $\mathcal{C}(2(k-1))$. Let $\mathcal{Aff}(\mathbb{R}^d, \mathbb{R})$ denote the space of affine maps between \mathbb{R}^d and \mathbb{R} .

It is worth pointing out that every map $\ell 1_C 1_{\mathcal{B}(0,1)}$ involved in the above decomposition of $R(\mathbf{c}, \mathbf{c}^*, \cdot)$ admits F_r as an envelop.

Denote by

$$\mathcal{G}'(r) = \left\{ \frac{\ell 1_C 1_{\mathcal{B}(0,M)}}{F_r}; \ell \in \mathcal{Aff}(\mathbb{R}^d, \mathbb{R}), C \in \mathcal{C}(2(k-1)) \right\}.$$

We immediately deduce that

$$\mathcal{N}(\mathcal{G}(r), \varepsilon) \leq \left(\mathcal{N}(\mathcal{G}'(r), \varepsilon/k^2) \right)^{k^2}.$$

Consider now the set of functions $\mathcal{N}(\mathcal{G}'(r), \varepsilon)$. The following lemma offers a decomposition of $\mathcal{N}(\mathcal{G}'(r), \varepsilon)$.

LEMMA 5.2. *Denote by \mathcal{F}_s $s = 1, \dots, p$ a collection of set of functions taking values in $[-1, 1]$. Then*

$$\mathcal{N} \left(\prod_{s=1}^p \mathcal{F}_s, \varepsilon \right) \leq \prod_{s=1}^p \mathcal{N}(\mathcal{F}_s, \varepsilon/p).$$

In order to apply Lemma 5.2, a crucial point is to only deal with maps taking values in $[-1, 1]$. To this aim, we define the set

$$\mathcal{H} = \left\{ \frac{f 1_{\{|f| \leq F_r\}}}{F_r} \right\},$$

where f is in $\mathcal{A}ff(\mathbb{R}^d, \mathbb{R})$, and the set

$$\mathcal{P}_d = \left\{ \left\{ 1_{\{\ell \leq a\}} \mid \ell \in \mathcal{L}(\mathbb{R}^d, \mathbb{R}), a \in \mathbb{R} \right\} \cup \left\{ 1_{\{\ell < a\}} \mid \ell \in \mathcal{L}(\mathbb{R}^d, \mathbb{R}), a \in \mathbb{R} \right\} \right\},$$

where $\mathcal{L}(\mathbb{R}^d, \mathbb{R})$ denotes the set of linear maps from \mathbb{R}^d to \mathbb{R} . We may write

$$\mathcal{G}'(r) \subset \mathcal{H} \times \prod_{i=1}^{2(k-1)} \mathcal{P}_d \times 1_{\mathcal{B}(0, M)},$$

It is well known that

$$d_p(\mathcal{P}_d) = d.$$

Since every set of functions in this decomposition is composed of functions taking values in $[-1, 1]$, we intend to apply Theorem 5.1 to every set. Consequently it remains to bound from above the pseudo-dimensions of these sets of functions.

First we deal with \mathcal{H} :

LEMMA 5.3. *One has*

$$d_p(\mathcal{H}) = d_p \left(\left\{ f 1_{\{|f| \leq F_r\}} \mid f \in \mathcal{A}ff(\mathbb{R}^d, \mathbb{R}) \right\} \right) \leq 24(3d + 4).$$

PROOF OF LEMMA 5.3. The first equality is obvious, so we only have to deal with the inequality. We recall that the pseudo-dimension of the set of functions $\{f 1_{\{|f| \leq F_r\}} \mid f \in \mathcal{A}ff(\mathbb{R}^d, \mathbb{R})\}$ is the Vapnik dimension of the set of functions

$$\mathcal{H}' = \left\{ 1_{\{f 1_{\{|f| \leq F_r\}} - t \leq 0\}} \mid f \in \mathcal{A}ff(\mathbb{R}^d, \mathbb{R}), t \in \mathbb{R} \right\}.$$

Let x_1, \dots, x_{2m} denote $2m$ points in \mathbb{R}^d . Since $F_r(x) = c_1 + c_2 1_{N^*(x)}$, where c_1 et c_2 are constants, at least m points fall in an area on which F_r takes the form $F_r(x) = c$, for some constant c . Without loss of generality, we assume that x_1, \dots, x_m fall in such an area. Consequently, we have to bound from above the quantity $\left| \left\{ 1_{\{f 1_{\{|f| \leq c\}} - t \leq 0\}} \right\} (x_1, \dots, x_m) \right|$. Observing that

$$\{1_{\{|f| \leq c\}}\} = \{1_{\{f \leq c\}} \times 1_{\{f \geq -c\}}\},$$

we deduce that

$$\begin{aligned} & \left| \left\{ 1_{\{f 1_{\{|f| \leq c\}} - t \leq 0\}} \right\} (x_1, \dots, x_m) \right| \\ & \leq \left| \left\{ 1_{\{f 1_{\{f \leq c\}} - t \leq 0\}} \right\} (x_1, \dots, x_m) \right| \times \left| \left\{ 1_{\{f 1_{\{f \geq -c\}} - t \leq 0\}} \right\} (x_1, \dots, x_m) \right|. \end{aligned}$$

Noticing that $d_{VC}(\{1_{\{f \leq c\}} | f \in \mathcal{A}ff(\mathbb{R}^d, \mathbb{R})\}) = d+1$ and making use of Sauer's lemma leads to, provided that $m \geq d+1$,

$$|\{1_{\{f \leq c\}}\}(x_1, \dots, x_m)| = |\{1_{\{f \geq -c\}}\}(x_1, \dots, x_m)| \leq \left(\frac{em}{d+1}\right)^{(d+1)},$$

which ensures that

$$|\{1_{\{|f| \leq c\}}\}(x_1, \dots, x_m)| \leq \left(\frac{em}{d+1}\right)^{2(d+1)}.$$

Choose a configuration of $\{1_{|f(x_1)| \leq c}, \dots, 1_{|f(x_m)| \leq c}\}$, for instance by indexing the x_i 's so that $|f(x_1)| > c, \dots, |f(x_r)| > c$ and $|f(x_{r+1})| \leq c, \dots, |f(x_m)| \leq c$. For the r first x_i 's, only two configurations remains for $\mathcal{H}'(x_1, \dots, x_r)$, the configuration $(0, \dots, 0)$ and $(1, \dots, 1)$. Considering the $m-r+1$ last x_i 's, $|\mathcal{H}'(x_{m-r+1}, \dots, x_m)| \leq |\{1_{\{f-t \leq 0\}}\}(x_{m-r+1}, \dots, m)|$. Next, $|\{1_{\{f-t \leq 0\}}\}(x_{m-r+1}, \dots, m)| \leq |\{1_{\{f-t \leq 0\}}\}(x_1, \dots, m)|$. Eventually, Sauer's lemma guarantees that $|\{1_{\{f-t \leq 0\}}\}(x_1, \dots, m)| \leq \left(\frac{em}{d+2}\right)^{(d+2)}$, provided that $m \geq d+2$. Consequently, we get, for $m \geq d+2$,

$$\begin{aligned} |\mathcal{H}'(x_1, \dots, x_{2m})| &= \left| \{1_{\{f 1_{\{|f| \leq F_r\}} - t \leq 0\}}\}(x_1, \dots, x_{2m}) \right| \\ &\leq 2^m \times \left| \{1_{\{f 1_{\{|f| \leq c\}} - t \leq 0\}}\}(x_1, \dots, x_m) \right| \\ &\leq 2^m \times |\{1_{\{|f| \leq c\}}\}(x_1, \dots, x_m)| \times 2 \left(\frac{em}{d+2}\right)^{d+2} \\ &\leq 2^m \times 2 \left(\frac{em}{d+1}\right)^{2(d+1)} \left(\frac{em}{d+2}\right)^{d+2}. \end{aligned}$$

To give an upper bound on $d_p(\mathcal{H})$, we have to find $m \geq d+2$ such that

$$2 \left(\frac{em}{d+1}\right)^{2(d+1)} \left(\frac{em}{d+2}\right)^{d+2} < 2^m.$$

Noticing that $x \mapsto \log_2(x)$ is a strictly concave map, we deduce that

$$2(d+1) \log_2 \left(\frac{em}{d+1}\right) + (d+2) \log_2 \left(\frac{em}{d+2}\right) < (3d+4) \log_2 \left(\frac{3em}{3d+4}\right).$$

Consequently, a sufficient condition on m is given by

$$\left(\frac{3em}{3d+4}\right)^{3d+4} \leq 2^{m-1}.$$

Using the same method as in [4], the choice $m = \lceil 3(3d+4) \log_2(3e) \rceil$ turns out to be adequate. At last, noticing that $2m \leq 24(3d+4)$, we immediately deduce that $d_p(\mathcal{H}) \leq 2m \leq 24(3d+4)$. \square

Applying Lemma 5.2 and Theorem 5.1 yields that

$$\begin{aligned} \mathcal{N}(\mathcal{G}'(r), \varepsilon) &\leq \mathcal{N}(\mathcal{H}, \varepsilon/(2k-1)) \times \mathcal{N}(\mathcal{P}_d, \varepsilon/(2k-1))^{2k-2} \\ &\leq \left(\frac{4k-2}{\varepsilon} \right)^{K[24(3d+4)+(d+1)(2k-2)]}. \end{aligned}$$

At last, the result of Proposition 5.1 is given by

$$\begin{aligned} \mathcal{N}(\mathcal{G}(r), \varepsilon) &\leq \mathcal{N}(\mathcal{G}'(r), \varepsilon/k^2)^{k^2} \\ &\leq \left(\frac{2(2k-1)k^2}{\varepsilon} \right)^{Kk^2[24(3d+4)+2(d+1)(k-1)]}. \end{aligned}$$

5.4. *Proof of Proposition 4.5.* The proof of Proposition 4.5 is based on elementary properties of distributions with finite support, which are extended to the case where the source distribution is supported on small balls. Throughout this subsection, a source distribution $P_{\sigma'}$ is fixed, so that $R(Q, P_{\sigma'})$ may be denoted by $R(Q)$.

LEMMA 5.4. *Let z_1 and z_2 be points in \mathbb{R}^d , denote by R the quantity $\|z_1 - z_2\|$, by U_i the ball $\mathcal{B}(z_i, \rho)$. At last, let P denote the cone-shaped distribution with density*

$$\frac{2(d+1)}{V} (\mathbb{K}_{\|x-z_i\| \leq \rho}(\rho - \|x - z_i\|)),$$

over each ball U_i , where V denote the volume of the unit ball. Then, if

$$(R/2 - 3\rho)^2 \geq \rho^2 \frac{2d(d+1)}{(d+2)(d+3)} \quad \text{and} \quad \rho \leq \frac{R}{2},$$

then the best 2-quantizer Q_2^ is such that $Q_2^*(U_i) = z_i$ for $i = 1, 2$. Furthermore, the best 1-quantizer Q_1^* is such that $Q_1^*(U_1 \cup U_2) = (z_1 + z_2)/2$.*

PROOF OF LEMMA 5.4. Let V_i denote the Voronoi cell associated with z_i in the Voronoi diagram generated by (z_1, z_2) . Denote by Q_2^* the quantizer satisfying $Q_2^*(U_i) = z_i$ for $i = 1, 2$.

For any quantizer Q denote by $R_i(Q) = \int_{V_i} \|x - Q(x)\|^2 dx$ the contribution of the cell i to the distortion of Q . Denote by V the volume of the

unit ball, and by S its surface. Recalling that $S = d \times V$, an elementary calculation shows that

$$\begin{aligned} R_i(Q_2^*) &= \frac{1}{2} \frac{d+1}{\rho^{d+1}V} \int_0^\rho S(\rho r^{d+1} - r^{d+2}) dr \\ &= \rho^2 \frac{d(d+1)}{2(d+2)(d+3)}. \end{aligned}$$

Let $m_i^{in} = |Q(U_i) \cap V_i|$ and $m_i^{out} = |Q(U_i) \cap V_i^c|$ denote the number of images of U_i sent inside and outside V_i . For a given i , there are three situations of interest, which are described below.

1. $m_i^{out} = 0$ and $m_i^{in} = 1$, then it is clear that $R_i(Q_2^*) \leq R_i(Q)$.
2. $m_i^{out} = 0$ and $m_i^{in} = 2$, then $R_i(Q) \geq 0 = R_i(Q_2^*) - \rho^2 \frac{d(d+1)}{2(d+2)(d+3)}$.
3. $m_i^{out} \geq 1$, then there exists $z \in U_i$ such that $Q(z) \notin V_i$. Consequently, $\|z - Q(z)\| \geq \frac{R}{2} - \rho$. Let $z' \in \mathcal{B}(z_i, \rho)$, then

$$\|z' - Q(z')\| \geq \|z - Q(z')\| - 2\rho \geq \|z - Q(z)\| - 2\rho \geq \frac{R}{2} - 3\rho.$$

Hence we deduce

$$R_i(Q) \geq 1/2 \left(\frac{R}{2} - 3\rho \right)^2 = R_i(Q_2^*) + 1/2 \left(\left(\frac{R}{2} - 3\rho \right)^2 - \rho^2 \frac{d(d+1)}{(d+2)(d+3)} \right).$$

Since Q is a 2-quantizer, it is easy to see that

$$|\{i; m_i^{in} \geq 2\}| \leq |\{i; m_i^{out} \geq 1\}|.$$

From this we deduce that

$$\begin{aligned} R(Q) &= \sum_{\{i; m_i^{in} \geq 2, m_i^{out} = 0\}} R_i(Q) + \sum_{\{i; m_i^{out} \geq 1\}} R_i(Q) + \sum_{\{i; m_i^{in} = 1, m_i^{out} = 0\}} R_i(Q) \\ &\geq R(Q_2^*) + \sum_{\{i; m_i^{in} \geq 2, m_i^{out} = 0\}} \frac{1}{2} \left(\left(\frac{R}{2} - 3\rho \right)^2 - \rho^2 \frac{d(d+1)}{(d+2)(d+3)} \right). \end{aligned}$$

Taking $(R/2 - 3\rho)^2 \geq \rho^2 \frac{2d(d+1)}{(d+2)(d+3)}$ ensures that $R(Q) \geq R(Q_2^*)$. \square

Considering the distributions P_σ , σ in $\{-1, +1\}^m$, taking $\rho \leq \frac{\Delta}{16}$ ensures that the conditions of Lemma 5.4 are satisfied when considering $P_{\sigma|U_i \cup U_i'}$. We turn now to the proof of Proposition 4.5.

Let Q be a k -quantizer. The following construction provides $Q_\sigma \in \mathcal{Q}$ such that $R(Q_\sigma) \leq R(Q)$. Let V_i denote the union of the Voronoi cells associated

with z_i and $z_i + \omega_i$, in the Voronoi diagram generated by the sequences z and ω . We adopt the following notation

$$\begin{cases} n_i(Q) &= |Q(\mathcal{B}(0, M)) \cap V_i| \\ n_i^{out}(Q) &= |Q(V_i) \cap V_i^c| \\ I_j(Q) &= \{i; n_i(Q) = j\} \\ i_j(Q) &= |I_j(Q)| \\ i_{\geq j}(Q) &= \sum_{i \geq j} i_j(Q). \end{cases}$$

The first step is to add code points to empty cells. From the k -quantizer Q , a quantizer Q_1 is built as follows

- If $n_i(Q) \geq 1$, then we take $Q_{1|V_i} \equiv Q|_{V_i}$.
- If $n_i(Q) = 0$, then we set $Q_1(U_i) = Q_1(U'_i) = z_i + \frac{w_i}{2}$.

Notice that Q_1 is a $(k + i_0(Q))$ -quantizer. Denote $k_1 = k + i_0$ and $p_{\pm} = \frac{1 \pm \delta}{2m}$, then $R(Q_1)$ can be bounded as follows.

Let i be an integer between 1 and m . We denote by $R_i(Q)$ the contribution of V_i to the risk $R(Q)$. If $i \in I_{\geq 1}$, then $R_i(Q) = R_i(Q_1)$. Otherwise, if $i \in I_0(Q)$,

$$R_i(Q_1) = 2p_{\pm}\rho^2 \frac{d(d+1)}{(d+2)(d+3)} + p_{\pm} \frac{\Delta^2}{2}.$$

Furthermore, if $i \in I_0$, then $n_i^{out}(Q) \geq 1$, which ensures that, as in the proof of Lemma 5.4,

$$R_i(Q) \geq p_{\pm} \left(\frac{(A-2)\Delta}{2} - 2\rho \right)^2.$$

Since $A \geq 6$ and $\rho \leq \frac{\Delta}{16}$, we may write

$$\begin{aligned} R_i(Q) - R_i(Q_1) &\geq p_{\pm} \left[(2\Delta - 2\rho)^2 - 4\rho^2 - \frac{\Delta^2}{2} \right] \\ &\geq p_{\pm} \left[2\Delta \left(\frac{3\Delta}{4} \right) - \frac{\Delta^2}{2} \right] \\ &\geq p_{-} \frac{3\Delta^2}{2}. \end{aligned}$$

Summing all the contributions of V_i 's leads to

$$R(Q_1) \leq R(Q) - i_0(Q)p_{-} \frac{3\Delta^2}{2}$$

Next, we build the quantizer Q_2 according to the following rule:

- If $n_i(Q_1) \geq 2$, then $Q_2(U_i) = z_i$ and $Q_2(U'_i) = z_i + w_i$.
- If $n_i(Q_1) = 1$, then $Q_2(U_i) = Q_2(U'_i) = z_i + \frac{w_i}{2}$.

Since for $i = 1, \dots, k$, $n_i(Q_1) \geq 1$, Q_2 is defined on every V_i . Notice that, since $I_j(Q_1) = I_j(Q)$ for $j \geq 2$, Q_2 has $k_2 = k + i_0(Q) - \sum_{p \geq 3} (p-2)i_p(Q)$ code points. The following lemma offers a relation between $R(Q_2)$ and $R(Q_1)$.

LEMMA 5.5. *One has*

$$R(Q_2) \leq R(Q_1) + i_{\geq 3}(Q) \frac{p_+ \Delta^2}{128}.$$

PROOF OF LEMMA 5.5. Let i be an integer between 1 and m . Several cases may occur, as described below.

- Assume that $n_i(Q_1) = 1$.
 - If $n_i^{out}(Q_1) = 0$, then $R_i(Q_1) \geq R_i(Q_2)$, according to Lemma 5.4.
 - If $n_i^{out}(Q_1) \geq 1$, then, using the same technique as mentioned to bound $R(Q_1)$ from above, $R_i(Q_1) - R_i(Q_2) \geq p_{\pm} \frac{3\Delta^2}{2}$, which leads to $R_i(Q_1) \geq R_i(Q_2)$.
- Assume that $n_i(Q_1) = 2$.
 - If $n_i^{out}(Q_1) = 0$, then $R_i(Q_1) \geq R_i(Q_2)$, according to Lemma 5.4.
 - If $n_i^{out}(Q_1) \geq 1$, then, since $R_i(Q_2) = 2p_{\pm} \frac{p^2 d}{d+2} \leq p_+ \frac{\Delta^2}{128}$, $R_i(Q_1) - R_i(Q_2) \geq \Delta^2 \geq 0$.
- At last, assume that $n_i(Q_i) \geq 3$. If $n_i^{out}(Q_1) \geq 1$, then $R_i(Q_1) \geq R_i(Q_2)$. If $n_i^{out}(Q_1) = 0$, then $R_i(Q_1) \geq 0 = R_i(Q_1) - 2p_{\pm} \frac{\Delta^2}{128}$. In both cases $R(Q_2) \leq R(Q_1) + p_+ \frac{\Delta^2}{128}$.

Noticing that $I_{\geq 3}(Q_1) = I_{\geq 3}(Q)$, and summing the contributions $R_i(Q_2)$ leads to the desired result. \square

The last step is to build a quantizer Q_{σ} from Q_2 with exactly k code points.

- If $k_2 = k$, set $Q_{\sigma} = Q_2$.
- If $k_2 < k$, choose $(k - k_2)$ V_i such that $n_i(Q_2) = 1$ (elementary calculation shows that there exist at least $k - k_2$ such V_i 's). For every such V_i , set $Q_{\sigma}(U_i) = z_i$ and $Q_{\sigma}(U'_i) = z_i + \omega_i$. Then

$$R(Q_{\sigma}) \leq R(Q_2) - (k - k_2)p_- \frac{\Delta^2}{2}.$$

- If $k_2 > k$, choose $(k_2 - k)$ cells V_i such that $n_i(Q_2) = 2$. For every such V_i , define $Q_\sigma(U_i) = Q_\sigma(U'_i) = z_i + \frac{\omega_i}{2}$. Then

$$R(Q_\sigma) \leq R(Q_2) + (k_2 - k)p_+ \frac{\Delta^2}{2}.$$

By construction, Q_σ has exactly k code points, and is an element of \mathcal{Q} . Finally, a result on the risk of Q_σ is given by the following proposition.

PROPOSITION 5.2. *Let Q be a quantizer and Q_σ built as mentioned above. Then,*

$$R(Q_\sigma) \leq R(Q).$$

PROOF OF PROPOSITION 5.2. Since $\delta \leq \frac{1}{3}$, easy calculation ensures that $1 - \frac{p_-}{p_+} \leq \frac{1}{2}$.

Suppose that $k_2 \leq k$. Comparing the risk of Q to the risks of Q_1 , Q_2 and Q_σ leads to

$$R(Q_\sigma) \leq R(Q) - i_0 p_- \frac{3\Delta^2}{2} + (i_0 + 2i_{\geq 3} - \sum_{p \geq 3} p i_p) p_- \frac{\Delta^2}{2} + i_{\geq 3} p_+ \frac{\Delta^2}{128}.$$

Since $\sum_{p \geq 3} p i_p \geq 3i_{\geq 3}$, it is clear that

$$\begin{aligned} R(Q_\sigma) &\leq R(Q) - p_- i_0 \frac{\Delta^2}{2} + \Delta^2 i_{\geq 3} \left(\frac{p_+}{128} - \frac{p_-}{2} \right) \\ &\leq R(Q). \end{aligned}$$

Next, suppose that $k_2 > k$. Then

$$\begin{aligned} R(Q_\sigma) &\leq R(Q) + \left(i_0 + 2i_{\geq 3} - \sum_{p \geq 3} p i_p \right) p_+ \frac{\Delta^2}{2} + i_{\geq 3} p_+ \frac{\Delta^2}{128} - i_0 p_- \frac{3\Delta^2}{2} \\ &\leq R(Q) + i_0 \frac{\Delta^2}{2} (p_+ - 3p_-) + p_+ i_{\geq 3} \Delta^2 \left(\frac{1}{128} - \frac{1}{2} \right), \end{aligned}$$

which yields that $R(Q_\sigma) \leq R(Q)$. \square

5.5. *Proof of Lemma 4.5.* Let introduce, for distributions P and Q with densities f and g the affinity

$$\alpha(P, Q) = \int \sqrt{fg},$$

so that $H^2(P, Q) = 2(1 - \alpha(P, Q))$. Elementary calculation shows that, if $\rho(\sigma, \sigma') = 4$, then

$$\alpha(P_\sigma, P_{\sigma'}) = 1 + \frac{2}{m} \left(\sqrt{1 - \delta^2} - 1 \right) \geq 1 - \frac{2\delta^2}{m}.$$

Hence we deduce

$$\begin{aligned} H^2(P_\sigma^{\otimes n}, P_{\sigma'}^{\otimes n}) &= 2(1 - \alpha(P_\sigma^{\otimes n}, P_{\sigma'}^{\otimes n})) \\ &= 2(1 - \alpha^n(P_\sigma, P_{\sigma'})) \\ &\leq \frac{4n\delta^2}{m}. \end{aligned}$$

Finally, since $\rho(\tau, \tau') = 2$ implies $\rho(\sigma(\tau), \sigma(\tau')) = 4$, for τ, τ' in $\{-1, +1\}^{\frac{m}{2}}$, the first part of Lemma 4.5 is proved.

Next, for simplicity assume that σ is such that $\sigma_1 = \dots \sigma_{\frac{m}{2}} = +1$ and $\sigma_{\frac{m}{2}+1} = \dots = \sigma_m = -1$. Let \mathcal{S}^- and \mathcal{S}^+ denote the set of mistakes of σ' , that is

$$\begin{cases} \mathcal{S}^- &= \{i = 1, \dots, \frac{m}{2} \mid \sigma'_i = -1\} \\ \mathcal{S}^+ &= \{i = \frac{m}{2} + 1, \dots, m \mid \sigma'_i = +1\}. \end{cases}$$

Finally let s^+ and s^- respectively denote $|\mathcal{S}^+|$ and $|\mathcal{S}^-|$. Since $\sum_{i=1}^m \sigma'_i = 0$, it is clear that $s^+ = s^- := s$.

As in Subsection 5.4, let $R_i(Q_{\sigma'})$ denote the contribution of $U_i \cup U'_i$ to the distortion. Then, for i in \mathcal{S}^- , elementary calculation shows that

$$R_i(Q_{\sigma'}) = R_i(Q_\sigma) + \frac{(1 + \delta)\Delta^2}{4m}.$$

Symmetrically, for i in \mathcal{S}^+ ,

$$R_i(Q_{\sigma'}) = R_i(Q_\sigma) - \frac{(1 - \delta)\Delta^2}{4m}.$$

Summing with respect to i and taking into account that $s^+ = s^- = s$ leads to

$$R(Q_{\sigma'}) = R(Q_\sigma) + s \frac{\Delta^2 \delta}{2m}.$$

Remarking that $s = \frac{\rho(\sigma, \sigma')}{4}$ concludes the proof of Lemma 4.5.

REFERENCES

- [1] ANTOS, A. (2005). Improved minimax bounds on the test and training distortion of empirically designed vector quantizers. *IEEE Trans. Inform. Theory* **51** 4022–4032. [MR2239018 \(2007b:94149\)](#)
- [2] ANTOS, A., GYÖRFI, L. and GYÖRGY, A. (2005). Individual convergence rates in empirical vector quantizer design. *IEEE Trans. Inform. Theory* **51** 4013–4022. [MR2239017 \(2007a:94125\)](#)
- [3] BARTLETT, P. L., LINDER, T. and LUGOSI, G. (1998). The minimax distortion redundancy in empirical quantizer design. *IEEE Trans. Inform. Theory* **44** 1802–1813. [MR1664098 \(2001f:94006\)](#)
- [4] BAUM, E. B. and HAUSSLER, D. (1989). What size net gives valid generalization? *Neural Comput.* **1** 151–160.
- [5] BIAU, G., DEVROYE, L. and LUGOSI, G. (2008). On the performance of clustering in Hilbert spaces. *IEEE Trans. Inform. Theory* **54** 781–790. [MR2444554 \(2009m:68221\)](#)
- [6] BLANCHARD, G., BOUSQUET, O. and MASSART, P. (2008). Statistical performance of support vector machines. *Ann. Statist.* **36** 489–531. [MR2396805 \(2009m:62085\)](#)
- [7] CHICHIGNOUD, M. and LOUSTAU, S. (2013-06). Adaptive Noisy Clustering.
- [8] FISCHER, A. (2010). Quantization and clustering with Bregman divergences. *J. Multivariate Anal.* **101** 2207–2221. [MR2671211 \(2012c:62188\)](#)
- [9] GERSHO, A. and GRAY, R. M. (1991). *Vector quantization and signal compression*. Kluwer Academic Publishers, Norwell, MA, USA.
- [10] GRAF, S. and LUSCHGY, H. (2000). *Foundations of quantization for probability distributions. Lecture Notes in Mathematics* **1730**. Springer-Verlag, Berlin. [MR1764176 \(2001m:60043\)](#)
- [11] KOLTCHINSKII, V. (2006). Local Rademacher complexities and oracle inequalities in risk minimization. *Ann. Statist.* **34** 2593–2656. [MR2329442 \(2009h:62060\)](#)
- [12] LEVRARD, C. (2013). Fast rates for empirical vector quantization. *Electron. J. Stat.* **7** 1716–1746.
- [13] LINDER, T. (2002). Learning-theoretic methods in vector quantization. In *Principles of nonparametric learning (Udine, 2001). CISM Courses and Lectures* **434** 163–210. Springer, Vienna. [MR1987659 \(2004f:68128\)](#)
- [14] LINDER, T., LUGOSI, G. and ZEGER, K. (1994). Rates of convergence in the source coding theorem, in empirical quantizer design, and in universal lossy source coding. *IEEE Trans. Inform. Theory* **40** 1728–1740. [MR1322387 \(96b:94005\)](#)
- [15] MAMMEN, E. and TSYBAKOV, A. B. (1999). Smooth discrimination analysis. *Ann. Statist.* **27** 1808–1829. [MR1765618 \(2001i:62074\)](#)
- [16] MASSART, P. (2007). *Concentration inequalities and model selection. Lecture Notes in Mathematics* **1896**. Springer, Berlin. Lectures from the 33rd Summer School on Probability Theory held in Saint-Flour, July 6–23, 2003, With a foreword by Jean Picard. [MR2319879 \(2010a:62008\)](#)
- [17] MASSART, P. and NÉDÉLEC, É. (2006). Risk bounds for statistical learning. *Ann. Statist.* **34** 2326–2366. [MR2291502 \(2009e:62282\)](#)
- [18] MENDELSON, S. and VERSHYNIN, R. (2003). Entropy and the combinatorial dimension. *Invent. Math.* **152** 37–55. [MR1965359 \(2004d:60047\)](#)
- [19] POLLARD, D. (1982). Quantization and the method of k -means. *IEEE Transactions on Information Theory* **28** 199–204.

- [20] POLLARD, D. (1982). A central limit theorem for k -means clustering. *Ann. Probab.* **10** 919–926. [MR672292 \(84c:60047\)](#)
- [21] TSYBAKOV, A. B. (2009). *Introduction to nonparametric estimation. Springer Series in Statistics*. Springer, New York. Revised and extended from the 2004 French original, Translated by Vladimir Zaiats. [MR2724359 \(2011g:62006\)](#)